

Core-Elements for Large-Scale Least Squares Estimation

Mengyu Li¹, Jun Yu^{*2}, Tao Li¹, and Cheng Meng^{†3}

¹Institute of Statistics and Big Data, Renmin University of China, Beijing, China

²School of Mathematics and Statistics, Beijing Institute of Technology, Beijing, China

³Center for Applied Statistics, Institute of Statistics and Big Data, Renmin University of China, Beijing, China

Abstract

The coresets approach, also called subsampling or subset selection, aims to select a subsample as a surrogate for the observed sample and has found extensive applications in large-scale data analysis. Existing coresets methods construct the subsample using a subset of rows from the predictor matrix. Such methods can be significantly inefficient when the predictor matrix is sparse or numerically sparse. To overcome this limitation, we develop a novel element-wise subset selection approach, called core-elements, for large-scale least squares estimation. We provide a deterministic algorithm to construct the core-elements estimator, only requiring an $O(\text{nnz}(X) + rp^2)$ computational cost, where X is an $n \times p$ predictor matrix, r is the number of elements selected from each column of X , and $\text{nnz}(\cdot)$ denotes the number of non-zero elements. Theoretically, we show that the proposed estimator is unbiased and approximately minimizes an upper bound of the estimation variance. We also provide an approximation guarantee by deriving a coresets-like finite sample bound for the proposed estimator. To handle potential outliers in the data, we further combine core-elements with the median-of-means procedure, resulting in an efficient and robust estimator with theoretical consistency guarantees. Numerical studies on various synthetic and real-world datasets demonstrate the proposed method's superior performance compared to mainstream competitors.

Keywords: Coresets, Linear model, Sparse matrix, Subset selection.

^{*}Joint first author

[†]Corresponding author, chengmeng@ruc.edu.cn

1 Introduction

Sparse matrices are matrices in which most of the elements are zero. Such matrices are common in various areas, including medical research, bioinformatics, privacy-preserving analysis, and distributed computing (Davis and Hu, 2011; Nguyen et al., 2023; Liu et al., 2021; Kairouz et al., 2021). In these areas, data are usually of high sparsity due to technical noises, data privacy concerns, and transmission cost, among others (Konečný et al., 2016; Andrews et al., 2021). One example is the single-cell RNA-sequencing (scRNA-seq) data containing information about the gene expression level of single cells. Owing to technical noises and intrinsic biological variability, scRNA-seq data expressed in count matrices always possess significant sparsity, known as the zero-inflation phenomenon (Nguyen et al., 2023). Another example is the word occurrence matrix, whose elements are calculated by multiplying two metrics, i.e., how many times a word appears in a document, and the inverse document frequency of the word across a set of documents. Such matrices are also highly sparse, especially for a short document and a large language model that contains millions of words (Qaiser and Ali, 2018).

In reality, many sparse matrices also exist due to missing values, i.e., the elements are not fully observed. Numerous methods have been developed to deal with the missing values in such cases, and most of these methods aim to fill the sparse matrix with some estimated non-zero values (Cai et al., 2010; Van Buuren and Groothuis-Oudshoorn, 2011; Hastie et al., 2015; Muzellec et al., 2020). In this paper, however, we are less concerned with missing values and instead focus on the case that the sparse matrix itself is fully observed.

We consider large-scale data analysis where the predictor matrix is highly sparse. One widely-used technique for large-scale data analysis is the coresets method, also called sub-sampling or subset selection. These methods select a subsample as a surrogate for the observed sample. Recently, such methods have been used pervasively in data reduction, measurement-constrained analysis, and active learning (Li and Meng, 2021; Meng et al., 2021; Settles, 2012). Various coresets methods have been proposed for linear regression (Dasgupta et al., 2009; Boutsidis et al., 2013; Ma and Sun, 2015; Meng et al., 2017; Dereziński et al., 2018; Ma et al., 2020; Wang et al., 2021), generalized linear regression (Wang et al., 2018; Ai et al., 2020, 2021; Yu et al., 2022), streaming time series (Xie et al.,

2019; Li et al., 2019), large-scale matrix approximation (Wang and Zhang, 2013; Alaoui and Mahoney, 2015; Wang et al., 2019), nonparametric regression (Ma et al., 2015; Meng et al., 2020; Sun et al., 2021; Meng et al., 2022; Dai et al., 2023), among others. Another avenue for handling large-scale data is data averaging (Wang et al., 2023), which is beyond the scope of this paper.

Despite the wide application, most existing coresets methods mainly focus on dense predictor matrices, and may be inefficient when the predictor matrix is of high sparsity. In particular, most of these methods construct the subsample using certain rows from the observed sample. When the observed predictor matrix is sparse, the selected subsample matrix also tends to be sparse for the subsampling methods that preserve the empirical distribution of the full data (Mak and Joseph, 2018; Joseph and Vakayil, 2022; Vakayil and Joseph, 2022). Such a subsample thus may lead to inefficient results, since the selected zero-valued elements have almost no impact on the down-streaming analysis, such as model estimation, prediction, and inference. Another category of subsampling approaches is designed for prediction purposes (Joseph and Mak, 2021; Chang, 2023; Dai et al., 2023), whose resulting subsample may not be statistically similar to the full data, thereby potentially overcoming the inefficiency discussed above. Nevertheless, this class of methods still encounters the probability of selecting a sparse subsample matrix when the full predictor matrix has an extremely high sparsity. Consequently, more efficient statistical tools suitable for sparse matrices are still meager.

In this paper, we bridge this gap by developing a novel element-wise subset selection method, called core-elements, for large-scale least squares estimation. Different from existing coresets methods that aim to select r rows from the predictor matrix $X \in \mathbb{R}^{n \times p}$ ($n \gg p$), we aim to construct a sparser subdata matrix $X^* \in \mathbb{R}^{n \times p}$ by keeping rp elements of X and zeroing out the remaining elements. Loosely speaking, our approach generalizes the existing coresets methods by getting rid of the requirement that the selected rp elements have to be located in r rows. Our major contributions are three-fold as follows.

- (1) We provide a deterministic algorithm to construct the core-elements estimator for linear regression. Utilizing such an estimator, we can approximate the least squares estimation within $O(\text{nnz}(X) + rp^2)$ computational time, where $\text{nnz}(\cdot)$ denotes the number of non-zero elements. Theoretical analysis demonstrates that our proposed estimator is

unbiased and approximately minimizes an upper bound on estimation variance.

- (2) We establish a coresets-like finite sample bound for the proposed estimator with approximation guarantees. In particular, we show that to achieve an $(1 + \epsilon)$ -relative error, the proposed estimator requires a subdata matrix X^* for the predictor matrix X such that the ratio $\|X - X^*\|_2 / \|X\|_2$ is $O(\epsilon^{1/2})$. Intuitively, such a result indicates that when X gets (numerically) sparser, fewer elements are required in X^* to achieve the $(1 + \epsilon)$ -relative error.
- (3) To handle potential outliers in the data, we develop a robust variant of the core-elements method by integrating it with the widely adopted median-of-means procedure (Lugosi and Mendelson, 2019; Lecué and Lerasle, 2020; Huang and Lederer, 2023). We further propose an algorithm to construct the robust estimator within $O(\text{nnz}(X) + rp^2)$ time. Theoretically, we show that the robust estimator is consistent with the true coefficient under certain regularity conditions.

To evaluate the empirical performance and computational efficiency of the proposed strategy, we compare it with mainstream competitors through extensive synthetic and real-world datasets, including uncorrupted and corrupted data with dense or sparse predictor matrices. Interestingly, although its primary design aims at sparse matrices, numerical findings reveal that the proposed estimator significantly outperforms its competitors regarding both estimation accuracy and CPU time, even when the predictor matrix is dense.

The remainder of this paper is organized as follows. We start in Section 2 by introducing the linear model and the state-of-the-art subsampling methods. In Section 3, we develop the core-elements estimator and present its theoretical properties. The robust version of the core-elements estimator is introduced in Section 4. We examine the performance of the proposed estimators through extensive simulations and a real-world example in Sections 5 and 6, respectively. Technical proofs, additional experimental results, and implementation code are provided in supplementary materials.

2 Background

We adopt the common convention of using uppercase letters for matrices and lowercase letters for vectors or scalars. We denote the Euclidean norm and ℓ_1 norm of a vector x as $\|x\|$ and $\|x\|_1$, respectively. For a matrix X , we represent its spectral norm (i.e., largest singular value) as $\|X\|_2$ and its Frobenius norm as $\|X\|_F$. The condition number of X , i.e., the ratio of its largest and smallest singular values, is denoted as $\kappa(X)$. Additionally, we use the notations $E(\cdot)$, $\text{pr}(\cdot)$, and $\text{tr}(\cdot)$ for mathematical expectation, probability measure, and trace, respectively.

2.1 Subsampling Methods for Least Squares Estimation

Consider the linear regression model,

$$y_i = x_i^\top \beta + \varepsilon_i, \quad i = 1, \dots, n. \quad (1)$$

Here $\{y_i\}_{i=1}^n \subset \mathbb{R}$ are the responses, $\{x_i\}_{i=1}^n \subset \mathbb{R}^p$ are the observations, $\beta \in \mathbb{R}^p$ is a vector of unknown coefficients, and $\{\varepsilon_i\}_{i=1}^n$ are independent and identically distributed (i.i.d.) error terms with zero mean and constant variance σ^2 . Let $y = (y_1, \dots, y_n) \in \mathbb{R}^n$ be the response vector, $X = (x_1, \dots, x_n)^\top \in \mathbb{R}^{n \times p}$ be the predictor matrix, and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$ be the noise vector. In this study, we assume that $n \gg p$, p is fixed, and X is of full column rank. We focus on the estimation of slope parameters and assume that the full data have been centralized. It is widely known that the ordinary least squares (OLS) estimator of β takes the form

$$\hat{\beta}_{\text{OLS}} = (X^\top X)^{-1} X^\top y. \quad (2)$$

In practice, the calculation of the least squares problem may suffer from high computational costs. Specifically, standard computation of the formulation (2) requires $O(np^2)$ computational time, which can be considerable when both n and p are large. To tackle the computational burden, various subsampling methods have been proposed. The main idea of subsampling methods can be described as follows: given a predictor matrix $X \in \mathbb{R}^{n \times p}$, subsample r rows (i.e., r observations) from X to construct a much smaller matrix $\hat{X} \in \mathbb{R}^{r \times p}$, and then use it as a surrogate for X in down-streaming analysis.

Most of the existing subsampling methods can be divided into two classes, i.e., the randomized subsampling approach and the design-based subsampling approach. The former class aims to carefully design a data-dependent non-uniform sampling probability distribution such that more informative data points will be selected with larger sampling weights (Ma et al., 2015; Meng et al., 2017; Knight, 2018; Ma et al., 2020; Ai et al., 2021). In contrast, the latter class aims to construct the most effective subsample estimator based on certain optimality criteria developed in the design of experiments (Wang et al., 2019; Meng et al., 2021; Wang et al., 2021, 2022; Chasiotis and Karlis, 2024). Despite the numerous subsampling algorithms proposed, most of them rely on row-wise sampling, which can be less effective when dealing with (numerically) sparse predictor matrices, as discussed in Section 1. In contrast, element-wise sampling is more adept at exploiting the inherent sparsity of data, motivating the development of the core-elements method. See Table 1 for a comprehensive comparison. We also refer readers to Li and Meng (2021) and Yu et al. (2023) for recent reviews.

Closely related to the subsampling methods are the coresets methods. These methods aim to select a subsample in a deterministic way, such that a loss function \mathcal{L} based on the subsample estimator is bounded by the loss function based on the full-sample estimator multiplying a constant $(1 + \epsilon)$ (Boutsidis et al., 2013; Munteanu et al., 2018; Feldman et al., 2020; Maalouf et al., 2022). In linear models, a subsample \hat{X} is called an $(1 + \epsilon)$ -coreset ($\epsilon > 0$), if there exists an estimator $\tilde{\beta}$ constructed by \hat{X} , such that

$$\|y - X\hat{\beta}_{\text{OLS}}\|^2 \leq \|y - X\tilde{\beta}\|^2 \leq (1 + \epsilon)\|y - X\hat{\beta}_{\text{OLS}}\|^2.$$

2.2 Sparse and Numerically Sparse Matrices

Recall that sparse matrices are matrices in which most elements are zero. Commonly, the sparsity is measured by using the ℓ_0 norm (i.e., the number of non-zero elements). Such a procedure, however, may not accurately reflect the simplicity of nearly sparse instances with a large number of small but non-zero elements. To combat the obstacle, existing literature also considers the so-called numerically sparse matrices (Gupta and Sidford, 2018; Braverman et al., 2021). Intuitively, numerically sparse is a weaker condition than sparse in which we do not require most elements to be zero, only that most elements are small

Table 1: Comparison of mainstream subsampling methods on computational cost, method type, and sampling type.

Method	Computational cost*	Method type	Sampling type
Full sample	$O(np^2)$	-	-
Uniform subsampling	$O(rp^2)$	Randomized	Row-wise
Doubly sketching (Hou-Liu and Browne, 2023)	$O(rp^2)$	Randomized	Row-wise
Leverage score subsampling (Ma and Sun, 2015)	$O(np^2 + rp^2)^\dagger$	Randomized	Row-wise
IBOSS (Wang et al., 2019)	$O(np + rp^2)$	Deterministic	Row-wise
OSS (Wang et al., 2021)	$O(np \log r + rp^2)$	Deterministic	Row-wise
D-optimal subsampling (Reuter and Schwabe, 2024)	$O(np^2 + rp^2)^\ddagger$	Deterministic	Row-wise
Core-elements (proposed)	$O(\text{nnz}(X) + rp^2)$	Deterministic	Element-wise

* Here, each method subsamples r rows or $s = rp$ elements from an $n \times p$ predictor matrix.

† The $O(np^2)$ component can be reduced to $O(np \log n)$ by involving some random projection-based approximation methods ([Drineas et al., 2012](#)).

‡ The $O(np^2)$ component can be reduced to $O(np)$ by using a simplified approximation that ignores non-diagonal elements of the covariance matrix.

enough to be ignored. Examples of numerically sparse are widely encountered in practice, including but not limited to linear programming constraints of the form $x_1 \geq \sum_{i=1}^n x_i/n$, and physical models whose strength of interaction decays with distance ([Carmon et al., 2020](#)).

2.3 Element-wise Sampling on Sparse Data Matrices

Consider the scenario that the observed predictor matrix $X \in \mathbb{R}^{n \times p}$ is a (numerically) sparse matrix in which most of the elements are (nearly) zero. Oftentimes, of interest is to find a sparser matrix $X^+ \in \mathbb{R}^{n \times p}$ that is a good proxy for X . Here, X^+ is also a sparse ma-

trix, from which the non-zero elements are a subset of the non-zero elements with respect to (w.r.t.) X . The problem of finding such a matrix X^+ has many applications in eigenvector approximation (Arora et al., 2006; Achlioptas and Mcsherry, 2007; El Karoui and d’Aspremont, 2010; Kundu et al., 2017; Gupta and Sidford, 2018), semi-definite programming (Arora et al., 2005; d’Aspremont, 2011; Garber and Hazan, 2016), matrix completion (Candès and Recht, 2009; Candès and Tao, 2010; Chen et al., 2014), optimal transport problems (Li et al., 2023a,b; Hu et al., 2024), and nonparametric regression (Li et al., 2024), among others.

To construct such a matrix X^+ , most existing methods aim to design a good sampling probability distribution p_{ij} ($i = 1, \dots, n; j = 1, \dots, p$), such that the approximation error $\|X - X^+\|_2$ is as small as possible given a sampling budget s . Achlioptas and Mcsherry (2007) proposed the so-called ℓ_2 sampling such that $p_{ij} \propto x_{ij}^2$, which was later refined by Drineas and Zouzias (2011). Alternatively, Arora et al. (2006) proposed ℓ_1 sampling, where $p_{ij} \propto |x_{ij}|$. Later, Achlioptas et al. (2013) proposed a near-optimal probability distribution using matrix-Bernstein inequality. Their approach can be regarded as a combination of two ℓ_1 -based distributions: $p_{ij} \propto |x_{ij}|$ when s is small, and $p_{ij} \propto |x_{ij}| \cdot \|x_i\|_1$ as s grows. More recently, Kundu et al. (2017) developed hybrid- (ℓ_1, ℓ_2) sampling, which is also a convex combination of probabilities previously proposed.

In general, existing literature on element-wise sampling mainly focus on the algorithmic perspective, such that the applications of interest are usually compressed sensing and matrix recovery. Nevertheless, element-wise sampling with a statistical perspective, such that the applications involve providing effective and efficient solutions for large-scale statistical modeling and inference, is still a blank field and remains further studied.

3 Core-Elements

In this section, we present our main algorithm and the theoretical properties of the proposed estimator. We first introduce the definition of core-elements and then construct an unbiased core-elements estimator. Following this, we derive an upper bound for the variance of such an estimator by utilizing the matrix-form Taylor expansion. Subsequently, we provide an algorithm for core-elements selection, resulting in an estimator that approx-

imately minimizes this upper bound. Furthermore, a coresets-like finite sample bound is provided to quantify the approximation error for the proposed estimator.

3.1 Problem Setup

We consider the problem of large-scale least squares estimation when the predictor matrix X is sparse or numerically sparse. To utilize the sparse structure, we propose to construct a sparser subdata matrix $X^* \in \mathbb{R}^{n \times p}$ from X carefully, and use X^* to construct an efficient least squares approximation. In particular, given a positive integer $s (< np)$, let S be an $n \times p$ matrix such that its elements involve s ones and $np - s$ zeros. The subdata matrix X^* then can be formulated as $X^* = S \odot X$, where \odot represents the Hadamard product, i.e., element-wise product. In other words, X^* is produced by keeping s elements of X and zeroing out the remaining elements. We then propose a general estimator that takes the form

$$\tilde{\beta}(D) = DX^{*\top}y, \quad (3)$$

where $D \in \mathbb{R}^{p \times p}$ is a scaling matrix to be determined. For simplicity, suppose that both $X^\top X$ and $X^{*\top} X$ are invertible in context; otherwise, their inverse operation should be replaced by a generalized inverse.

A natural question arises: given a fixed budget s , among all the estimators that take the form (3), how to construct the scaling matrix D and the sparse subdata matrix X^* , such that (i) the estimator $\tilde{\beta}(D)$ is unbiased and (ii) its estimation variance is as small as possible? For the first part of this question, one can see that when $D = (X^{*\top} X)^{-1}$, $E\{\tilde{\beta}(D)\} = (X^{*\top} X)^{-1} X^{*\top} \{X\beta + E(\varepsilon)\} = \beta$, indicating that $\tilde{\beta}(D)$ is unbiased to β . Such a finding motivates us to focus on the unbiased estimator

$$\tilde{\beta} = (X^{*\top} X)^{-1} X^{*\top} y.$$

Consider the classical subsample-based estimator

$$\tilde{\beta}' = (\hat{X}^\top \hat{X})^{-1} \hat{X}^\top \hat{y},$$

where $\hat{X} \in \mathbb{R}^{r \times p}$ represents a subsample (of full column rank) consisting of r rows from X , and $\hat{y} \in \mathbb{R}^r$ contains the corresponding elements from y . When the selected $s = rp$

elements are located in r rows, i.e., when the element-wise subset selection degenerates to the row selection, it can be shown that $\tilde{\beta} = \tilde{\beta}'$. While in general cases, these two estimators are different.

Next, we consider the variance of the estimator $\tilde{\beta}$. Recall that σ^2 represents the variance of random errors in model (1). Our goal is to find a subdata matrix X^* that minimizes the estimation variance, which has a closed form

$$\begin{aligned} E(\|\tilde{\beta} - \beta\|^2|X) &= E\{y^\top X^*(X^\top X^*)^{-1}(X^{*\top} X)^{-1} X^{*\top} y|X\} - \beta^\top \beta \\ &= \beta^\top \beta + \sigma^2 \text{tr}\{X^*(X^\top X^*)^{-1}(X^{*\top} X)^{-1} X^{*\top}\} - \beta^\top \beta \\ &= \sigma^2 \|(X^{*\top} X)^{-1} X^{*\top}\|_F^2. \end{aligned} \quad (4)$$

However, this variance term is challenging to minimize directly. To overcome the obstacle, we provide an upper bound for (4) and aim to minimize the upper bound instead. The upper bound is derived by utilizing the matrix-form Taylor expansion (see Chapter 1, Higham (2008)), detailed in Lemma 1.

Lemma 1. *Let $L = X - X^*$. A Taylor expansion of $E(\|\tilde{\beta} - \beta\|^2|X)$ around the point $X^* = X$ yields the following upper bound,*

$$E(\|\tilde{\beta} - \beta\|^2|X) \leq \sigma^2[\text{tr}\{(X^\top X)^{-1}\} + \|(X^\top X)^{-1}\|_2^2 \|L\|_F^2] \{1 + O(\lambda_0)\}. \quad (5)$$

Here, assume the spectral radius $\lambda_0 = \|(X^\top X)^{-1} L^\top X\|_2 < 1$ to ensure the convergence of the matrix series.

When the Taylor expansion in Lemma 1 is valid, the inequality (5) indicates that the upper bound of the estimation variance decreases as $\|L\|_F$ and the remainder λ_0 decreases. Considering λ_0 , we have

$$\lambda_0 \leq \|(X^\top X)^{-1}\|_2 \|X\|_2 \|L\|_2 \leq \|(X^\top X)^{-1}\|_2 \|X\|_2 \left(p \max_{j \in \{1, \dots, p\}} l^{(j)\top} l^{(j)} \right)^{1/2},$$

where $l^{(j)}$ denotes the j th column of L . Such an inequality indicates that a smaller value of the maximum column norm of L is associated with a smaller λ_0 . As a result, to minimize the upper bound in Lemma 1, we need to keep both $\|L\|_F$ and the column norms of L as small as possible.

Motivated by this, we propose to construct X^* by keeping the elements with the top largest absolute values w.r.t. each column of X and zeroing out the remaining. Intuitively, for a fixed number of selected elements, such L has the approximately minimum column norm respecting every column. Thus, both the values of $\|L\|_F$ and $\|L\|_2$ will be approximately minimized, resulting in a relatively small upper bound of the estimation variance in Lemma 1. Moreover, the column-wise process can prevent any entire column of X from being discarded, thus avoiding producing a singular X^* and ensuring the estimability of coefficients. We call such a procedure “core-elements” since the idea behind it is analogous to “coresets”, except that what we select are elements instead of rows.

3.2 Main Algorithm

Without loss of generality, we assume the number of selected elements $s = rp$, where r is an integer. Algorithm 1 summarizes the construction of core-elements and the proposed estimator, which are illustrated in Fig. 1.

Algorithm 1 CORE-ELEMENTS(X, y, r)

- 1: **Input:** $X = (x_{ij}) \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$, $r \in \mathbb{Z}_+$
 - 2: Initialize $S = (0) \in \mathbb{R}^{n \times p}$
 - 3: **for** $j = 1, \dots, p$ **do**
 - 4: Let $\mathcal{I} = \{i_1, \dots, i_r\}$ be an index set, such that $\{|x_{i_q j}|\}_{q=1}^r$ are the r largest ones among $\{|x_{ij}|\}_{i=1}^n$
 - 5: Let $s_{i_q j} = 1$, $q = 1, \dots, r$
 - 6: **end for**
 - 7: $X^* = S \odot X$, where \odot represents the element-wise product
 - 8: **Return** $\tilde{\beta} = (X^{*\top} X)^{-1} X^{*\top} y$
-

Consider the computational cost of Algorithm 1. Constructing the matrix X^* using a partition-based selection algorithm requires $O(\text{nnz}(X))$ time (Musser, 1997; Martinez, 2004; Wang et al., 2019). Each column of X^* contains at most r non-zero elements, thus calculating $X^{*\top} X$ takes $O(rp^2)$ time by using sparse matrix representations and operations. The computing time for $\tilde{\beta}$ is thus at the order of $O(rp^2 + p^3)$. Therefore, the overall computational cost of Algorithm 1 is $O(\text{nnz}(X) + rp^2)$, which becomes $O(\text{nnz}(X))$ when

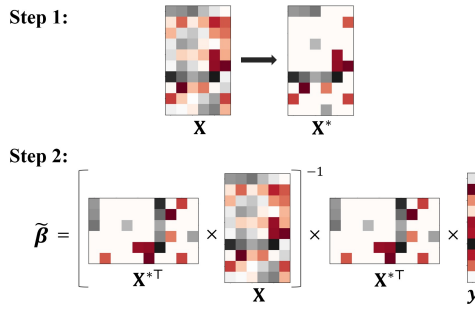


Figure 1: Illustration for Algorithm 1. Each element of a matrix is labeled with a different color, such that larger positive and smaller negative values are labeled with more red and more black colors, respectively. Step 1 illustrates the core-elements selection, where $r = 3$ elements with the largest absolute values w.r.t. each column are selected. Step 2 illustrates the proposed estimator.

$n \gg r$. In the large-sample scenario such that $n \gg p$, such an algorithm is much faster than the popular leverage-based subsampling methods (Ma et al., 2015). This is because the leverage-based methods involve the singular value decomposition of the full predictor matrix, requiring a cost of the order $O(np^2)$.

3.3 Theoretical Properties

We now present our main theorem, indicating that the selected elements are $(1 + \epsilon)$ -core-elements for least squares estimation. In other words, under some regularity conditions, the proposed estimate achieves the $(1 + \epsilon)$ -relative error w.r.t. the ℓ_2 loss. Technical proofs are provided in the supplementary material.

Theorem 2. *Given $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$ and $\epsilon > 0$, let X^* be the subdata matrix and $\tilde{\beta}$ be the estimate that calculated by Algorithm 1. When X^* satisfies $\|X - X^*\|_2 \leq \epsilon' \|X\|_2$ with*

$$0 < \epsilon' \leq \frac{1}{\kappa^2(X)} \left[1 + \frac{\{\kappa^2(X) + 1\} \|y\|}{\epsilon^{1/2} \|y - X\hat{\beta}_{OLS}\|} \right]^{-1}, \quad (6)$$

we have

$$\|y - X\hat{\beta}_{OLS}\|^2 \leq \|y - X\tilde{\beta}\|^2 \leq (1 + \epsilon) \|y - X\hat{\beta}_{OLS}\|^2. \quad (7)$$

Theorem 2 indicates that to achieve the $(1 + \epsilon)$ -relative error in inequality (7), Algorithm 1 requires a subdata matrix X^* such that the ratio $\|X - X^*\|_2 / \|X\|_2$ is $O(\epsilon^{1/2})$. This

rate is supported by empirical results in Section 5. Intuitively, such a result also indicates that when the predictor matrix X gets (numerically) sparser, fewer elements are required in X^* to achieve the same relative error w.r.t. the ℓ_2 loss.

In addition, the value of ϵ' also depends on the condition number $\kappa(X)$ and the relative sum of squares error (SSE) $\|y - X\hat{\beta}_{\text{OLS}}\|^2/\|y\|^2$. Specifically, a larger ϵ' is admitted to achieve the $(1 + \epsilon)$ -relative error if the condition number $\kappa(X)$ decreases and the SSE increases. The following remark discusses the relationship between r and ϵ' .

Remark 3. Suppose X is a sparse covariate matrix with $\alpha \times 100\%$ non-zero elements ($0 < \alpha \leq 1$), and each column has the same number of non-zero elements. Further, suppose the non-zero elements of X are i.i.d. from a continuous cumulative distribution function F . We consider two specific cases of the distribution F as follows.

- Case 1: uniform distribution over $(-1, 1)$.

If the subsample parameter r in Algorithm 1 satisfies $r < \alpha n$ and

$$\frac{r}{n} \geq \alpha - \frac{(\alpha\epsilon'\|X\|_2)^{2/3}}{(2np)^{1/3}},$$

then the subdata matrix X^* achieves the condition $\|X - X^*\|_2 \leq \epsilon'\|X\|_2$ in Theorem 2 with high probability for a relatively large n . Under a general condition $\|X\|_2 = O((np)^{1/2})$, which can be satisfied as long as all the elements in X are bounded, such a result indicates that Algorithm 1 needs to select around $\{\alpha - (c\alpha\epsilon')^{2/3}\} \times 100\%$ elements to achieve the $(1 + \epsilon)$ -relative error for a constant $c > 0$.

- Case 2: standard normal distribution on \mathbb{R} .

If $r < \alpha n$ satisfies

$$\frac{r}{n} \geq \alpha - \min \left\{ \alpha\phi, \frac{(\epsilon'\|X\|_2)^2}{2G^{-1}(\phi)np} \right\}, \quad (8)$$

where $0 < \phi < 1$ and G is the cumulative distribution function of the chi-squared distribution with 1 degree of freedom, then the condition $\|X - X^*\|_2 \leq \epsilon'\|X\|_2$ holds with high probability when n is sufficiently large. To achieve the $(1 + \epsilon)$ -relative error, (8) indicates us to select $[\alpha - \min\{\alpha\phi, (c\epsilon')^2/G^{-1}(\phi)\}] \times 100\%$ elements under the condition $\|X\|_2 = O((np)^{1/2})$, where $c > 0$ is a constant.

4 MOM Core-Elements

To guarantee the robustness of the core-elements approach, we modify Algorithm 1 by combining it with the popular median-of-means (MOM) procedure (Hsu and Sabato, 2014; Lugosi and Mendelson, 2019; Lecué and Lerasle, 2019, 2020; Mathieu, 2021; Huang and Lederer, 2023) to provide a robust version of core-elements, called “MOM core-elements”. We first propose an algorithm in the divide-and-conquer framework to obtain the MOM core-elements estimator, and then we establish consistency of the proposed estimator under regularity conditions.

4.1 Proposed Algorithm

Due to the existence of outliers, we relax the standard i.i.d. setup to the $\mathcal{I} \cup \mathcal{O}$ framework following the work of Lecué and Lerasle (2020). Specifically, we assume that data are partitioned into two (unknown) sets, \mathcal{I} and \mathcal{O} , such that $\mathcal{I} \cup \mathcal{O} = \{1, \dots, n\}$ and $\mathcal{I} \cap \mathcal{O} = \emptyset$. Data $\{(x_i, y_i)\}_{i \in \mathcal{I}}$ are i.i.d. informative data, and data $\{(x_i, y_i)\}_{i \in \mathcal{O}}$ are outliers on which no assumption is granted.

Given a random partition of the full data $(X, y) = \{(x_i, y_i)\}_{i=1}^n$ into blocks of equal sizes, the principle of MOM estimator is that we first obtain an estimation on each block independently, and then we aggregate the results from these blocks by taking the median. Let k be the number of blocks and n_l ($l = 1, \dots, k$) be the size of each block. Without loss of generality, we assume that both n and r are divisible by k . Then $n_l \equiv n/k$ for $l = 1, \dots, k$. Denote $(X^{(l)}, y^{(l)}) \in \mathbb{R}^{n_l \times p} \times \mathbb{R}^{n_l}$ be the data partitioned into the l th block. For the l th block, we construct the core-elements matrix $X^{(l)*}$ containing $r_l p$ non-zero elements, where $r_l \equiv r/k$ for $l = 1, \dots, k$. Next, we obtain the corresponding estimator $\tilde{\beta}^{(l)} = (X^{(l)*\top} X^{(l)})^{-1} X^{(l)*\top} y^{(l)}$. Then, the MOM core-elements estimator is defined as

$$\tilde{\beta}_{\text{MOM}} = \text{med}(\tilde{\beta}^{(1)}, \dots, \tilde{\beta}^{(k)}), \quad (9)$$

where $\text{med}(\cdot)$ denotes the coordinate-wise median. Algorithm 2 details the MOM core-elements procedure.

In Algorithm 2, the number of blocks k indicates the robustness of our proposed estimator. A popular measure to quantify robustness is the breakdown point (Hampel, 1968;

Algorithm 2 MOM CORE-ELEMENTS(X, y, r, k)

- 1: **Input:** $X = (x_{ij}) \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$, $r, k \in \mathbb{Z}_+$
 - 2: Partition (X, y) into k blocks $\{(X^{(l)}, y^{(l)})\}_{l=1}^k$ randomly and evenly
 - 3: **for** $l = 1, \dots, k$ **do**
 - 4: Compute $\tilde{\beta}^{(l)} = \text{CORE-ELEMENTS}(X^{(l)}, y^{(l)}, r/k)$
 - 5: **end for**
 - 6: **Return** $\tilde{\beta}_{\text{MOM}} = \text{med}(\tilde{\beta}^{(1)}, \dots, \tilde{\beta}^{(k)})$
-

Donoho and Huber, 1983; Donoho and Gasko, 1992; Lecué and Lerasle, 2020), defined as the smallest proportion of corrupted observations needed to push an estimator to infinity. The breakdown point of Algorithm 2 is $\lfloor k/2 \rfloor / n$, because less than $\lfloor k/2 \rfloor$ outliers may corrupt at most $\lfloor k/2 \rfloor$ blocks, leaving the median in (9) equal to the estimation on a single block with uncorrupted data. Remark that Algorithm 1 is a special case of Algorithm 2 when $k = 1$, which is applicable to datasets without extreme outliers.

Consider the computation time of Algorithm 2. Constructing the core-elements estimator on the l th block requires $O(\text{nnz}(X^{(l)}) + rp^2/k + p^3)$ time, and the aggregation step needs $O(kp)$ time. Thus, the total computational cost of Algorithm 2 is at the order of $O(\text{nnz}(X) + rp^2 + kp^3 + kp) = O(\text{nnz}(X) + rp^2)$, where the equation holds because r/k should be at least of the same order as p for ensuring the well-definedness of $\tilde{\beta}^{(1)}, \dots, \tilde{\beta}^{(k)}$. Such cost is the same as that of Algorithm 1.

4.2 Theoretical Properties

Now we provide the convergence of the MOM core-elements estimator. To begin with, we introduce the following regularity conditions.

- (H1) Denote the Fisher information matrix $I^{(l)} = n_l^{-1} X^{(l)\top} X^{(l)}$ on the l th block for $l = 1, \dots, k$. Assume that $c < \inf_l \lambda_{\min}(I^{(l)}) \leq \sup_l \lambda_{\max}(I^{(l)}) < \infty$ for some constant $c > 0$, where $\lambda_{\max}(\cdot)$ and $\lambda_{\min}(\cdot)$ respectively stand for the maximum and minimum eigenvalues.
- (H2) Let $L^{(l)} = X^{(l)} - X^{(l)*}$ for $l = 1, \dots, k$. Assume that $\sup_l \|L^{(l)}\|_F^2 / n_l^2 \rightarrow 0$ as $n_l \rightarrow \infty$.
- (H3) Assume that the spectral radius $\lambda_0^{(l)} = \|(X^{(l)\top} X^{(l)})^{-1} L^{(l)\top} X^{(l)}\|_2 < 1$ for $l = 1, \dots, k$.

(H4) Suppose that $k > 2|\mathcal{B}_O| + 1$, where \mathcal{B}_O is the set of blocks containing at least one outlier and $|\cdot|$ denotes the cardinal number. Further, assume that \mathcal{B}_O contains finite elements.

Theorem 4. *Suppose the conditions (H1)–(H4) hold almost surely. As $k, n_l \rightarrow \infty$ ($l = 1, \dots, k$), $\tilde{\beta}_{\text{MOM}}$ converges to β in probability, i.e., for any given $\epsilon > 0$, it holds that*

$$\text{pr}(\|\tilde{\beta}_{\text{MOM}} - \beta\| > \epsilon) \rightarrow 0.$$

In Theorem 4, condition (H1) is commonly assumed in the literature. Conditions (H2) and (H3) bound the error of the local core-elements estimator on each block by using Lemma 1. Condition (H4) guarantees the effectiveness of the MOM procedure according to its breakdown point discussed above.

5 Simulation Studies

In this section, we first evaluate the performance of core-elements (i.e., Algorithm 1) in estimating β and predicting y on uncorrupted synthetic datasets. Subsequently, we provide empirical evidence to support the error bound in Theorem 2. Next, we consider corrupted datasets to show the effectiveness and robustness of MOM core-elements (i.e., Algorithm 2). Finally, we demonstrate the advantage of the proposed strategy over other subsampling methods w.r.t. computational efficiency.

5.1 Performance on Uncorrupted Data

We use CORE to refer to the estimator in Algorithm 1. For comparison, we consider the full sample OLS estimation (FULL) and several state-of-the-art subsampling methods mentioned in Table 1, including uniform subsampling (UNIF), doubly sketching (DOUBLY) (Hou-Liu and Browne, 2023), basic leverage subsampling (BLEV) (Drineas et al., 2006; Ma et al., 2015), shrinkage leverage subsampling (SLEV) with shrinkage parameter being 0.9 (Ma et al., 2015), information-based optimal subset selection (IBOSS) (Wang et al., 2019), orthogonal subsampling (OSS) (Wang et al., 2021), and D-optimal subsampling (DOPT) (Reuter and Schwabe, 2024).

For uncorrupted data, the predictor matrix X is generated from different kinds of widely-used distributions:

(D1) multivariate normal distribution, $N(0_p, \Sigma)$;

(D2) multivariate log-normal distribution, $LN(0_p, \Sigma)$;

(D3) multivariate t-distribution with 3 degrees of freedom, $t_3(0_p, \Sigma)$,

where $\Sigma = (\sigma_{ij}) \in \mathbb{R}^{p \times p}$ is a covariance matrix with $\sigma_{ij} = 0.6^{|i-j|}$ for $i, j = 1, \dots, p$.

To introduce sparsity, after generating the predictor matrix and centering it, we randomly zero out their elements with a sparsity ratio α . Specifically, we randomly select $\alpha \times 100\%$ of the elements and set these elements to be zero. Consider $\alpha = \{0, 0.2, 0.4, 0.6, 0.8\}$, referred to as (R1)–(R5), respectively. Next, we add a small random perturbation following $U(-0.1, 0.1)$ to each zero element of the predictor matrix to obtain a numerically sparse matrix. Then, (R1) corresponds to a completely dense matrix, and (R5) corresponds to a highly numerically sparse matrix. We then generate the response y from the linear model (1). The true coefficient β is a p -dimensional vector of ones, and the signal-to-noise ratio, defined as $\text{SNR} = \text{var}(X\beta)/\sigma^2$, is set to be 4. The simulations in misspecified linear models and alternative choices of true β are relegated to the supplementary material. Let the sample size $n = 10^4$ and dimension $p = 10^2$. For the row-wise subsampling methods, we select $r \in \{2p, 4p, 6p, 8p, 10p\}$ rows for each of these methods. For a fair comparison, we select $s = rp$ elements for the proposed core-elements method.

We calculate the empirical mean squared error (MSE) for each of the estimators based on one hundred replications, i.e.,

$$\text{MSE} = \frac{1}{100} \sum_{i=1}^{100} \frac{\|\widehat{\beta}^{(i)} - \beta\|^2}{\|\beta\|^2}, \quad (10)$$

where $\widehat{\beta}^{(i)}$ represents the estimator in the i th replication. We also consider the prediction MSE (PMSE). For the i th replication, we randomly split the observed sample into a training set $(y_{\text{train}}^{(i)}, X_{\text{train}}^{(i)})$ of size $\lfloor 0.7n \rfloor$ and a test set $(y_{\text{test}}^{(i)}, X_{\text{test}}^{(i)})$ of size $\lceil 0.3n \rceil$. For each subsampling method, we select a subset from the training set leading to an estimator $\widehat{\beta}_{\text{train}}^{(i)}$, and then use it to predict the response y_{test} in the test set. In this way, PMSE is calculated

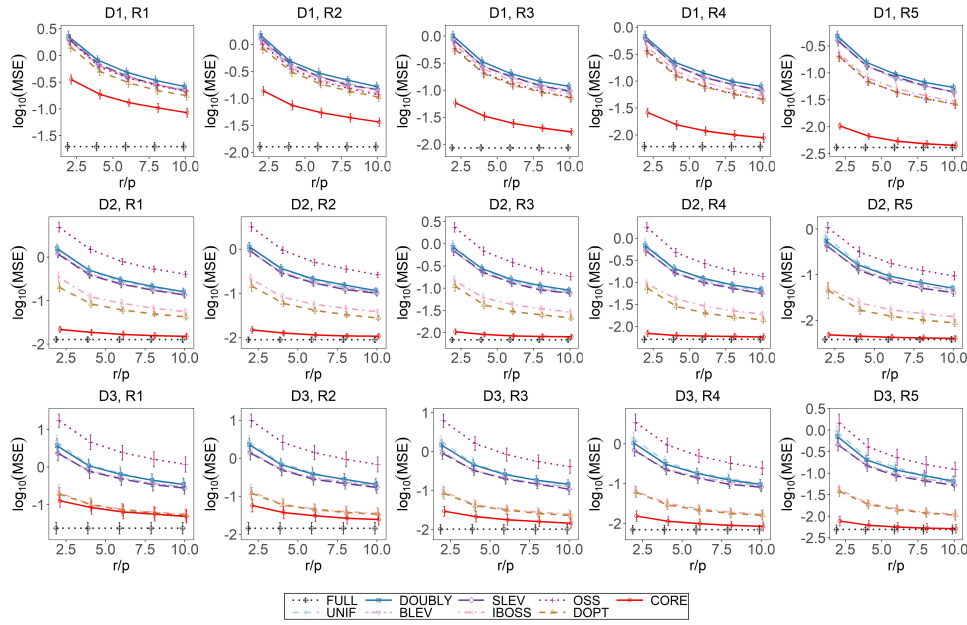


Figure 2: Comparison of different estimators w.r.t. MSE. Each row represents a particular data distribution, i.e., (D1)–(D3), and each column represents a different sparsity ratio, i.e., (R1)–(R5). Vertical bars are the standard errors.

as

$$\text{PMSE} = \frac{1}{100} \sum_{i=1}^{100} \frac{\|X_{\text{test}}^{(i)} \hat{\beta}_{\text{train}}^{(i)} - y_{\text{test}}^{(i)}\|^2}{\|y_{\text{test}}^{(i)}\|^2}. \quad (11)$$

The results of $\log(\text{MSE})$ and $\log(\text{PMSE})$ versus different subsample sizes are shown in Figs. 2 and 3, respectively.

In Figs. 2 and 3, we observe that both MSE and PMSE w.r.t. all estimators decreases as r increases. We also observe that CORE consistently outperforms other subsampling approaches under all circumstances. The advantage becomes more apparent when the level of sparsity increases, i.e., from (R1) to (R5). Such an observation indicates the proposed estimator can provide a more effective estimate than the competitors, especially when the predictor matrix is of high sparsity. Such success can be attributed to the fact that the proposed core-elements approach can effectively utilize the sparsity structure of the predictor matrix, and the proposed estimator is unbiased and has an approximately minimized estimation variance.

While other deterministic subsampling methods (i.e., IBOSS, OSS, and DOPT) are also competitive, their performance varies significantly across different data distributions.

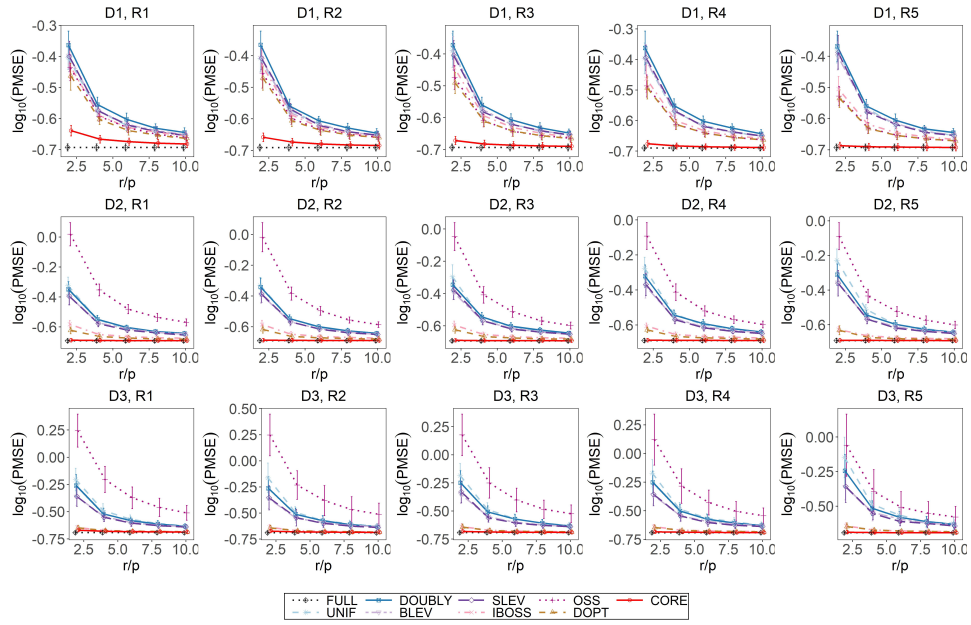


Figure 3: Comparison of different estimators w.r.t. PMSE. Each row represents a particular data distribution, i.e., (D1)–(D3), and each column represents a different sparsity ratio, i.e., (R1)–(R5). Vertical bars are the standard errors.

Specifically, the OSS method performs well with normally distributed data (D1), but it loses efficacy with asymmetric log-normal (D2) and heavy-tailed t-distributions (D3), performing worse than uniform sampling. This behavior aligns with the findings presented in Table 3 of Yu et al. (2023). In contrast, while IBOSS and DOPT are comparable to CORE in handling heavy-tailed distributions, they fall short in normal distributions. Overall, the proposed core-elements approach demonstrates both generality and superiority across a variety of data distributions.

In order to verify the theoretical error bound provided in Theorem 2, we compare the empirical and theoretical values of ϵ under different values of ϵ' , as shown in Fig. 4. Specifically, given a small ϵ' , the empirical value of ϵ is calculated as

$$\frac{\|y - X\tilde{\beta}\|^2}{\|y - X\hat{\beta}_{\text{OLS}}\|^2} - 1$$

according to (7), and the theoretical value of ϵ is calculated as

$$\left[\frac{\epsilon' \kappa^2(X) \{\kappa^2(X) + 1\} \|y\|}{\{1 - \epsilon' \kappa^2(X)\} \|y - X\hat{\beta}_{\text{OLS}}\|} \right]^2$$

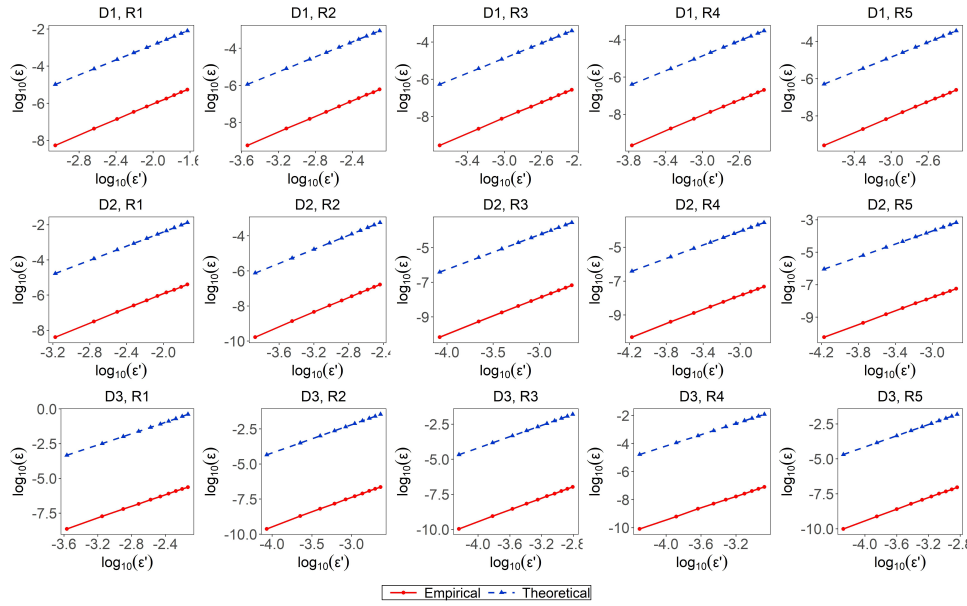


Figure 4: Comparison of the empirical value and the theoretical value of the error term ϵ . Each row represents a particular data distribution, i.e., (D1)–(D3), and each column represents a different sparsity ratio, i.e., (R1)–(R5).

according to (6). Both ϵ and ϵ' have been made logarithmic transformation in Fig. 4. We can observe that although the empirical and theoretical values of ϵ differ, their growth trends have an apparent parallel pattern. This observation indicates that our proposed error bound and the empirical value are of the same order. Their difference is up to a constant under the log transformation.

5.2 Performance on Corrupted Data

We compare the proposed MOM core-elements approach (i.e., Algorithm 2), referred to as MOM-CORE, with the full sample OLS estimation and the subsampling methods mentioned above in the presence of outliers. To ensure fairness, all competing methods are equipped with the MOM procedure.

The corrupted data consist of informative data $\{(x_i, y_i)\}_{i \in \mathcal{I}}$ and various types of outliers $\{(x_i, y_i)\}_{i \in \mathcal{O}}$ with $\mathcal{O} = \mathcal{O}_1 \cup \mathcal{O}_2 \cup \mathcal{O}_3 \cup \mathcal{O}_4$, such that $|\mathcal{O}| = n_o$ and $|\mathcal{I}| = n - n_o$. $\{(x_i, y_i)\}_{i \in \mathcal{I}}$ is generated in the same way as in the previous section, and $\{(x_i, y_i)\}_{i \in \mathcal{O}}$ is constructed following the setup in Lecué and Lerasle (2020). More precisely,

- for $i \in \mathcal{O}_1$ with $|\mathcal{O}_1| = \lceil n_o/4 \rceil$, $y_i = 1000 + 10\zeta_1$ and $x_i = -10 \times 1_p + \zeta_2$, where $\zeta_1 \in \mathbb{R}$ and $\zeta_2 \in \mathbb{R}^p$ are noises following the (multivariate) standard normal distribution;
- for $i \in \mathcal{O}_2$ with $|\mathcal{O}_2| = \lceil n_o/4 \rceil$, $y_i = -500 + 10\zeta_1$ and $x_i = 10 \times 1_p + \zeta_2$;
- for $i \in \mathcal{O}_3$ with $|\mathcal{O}_3| = \lceil n_o/4 \rceil$, y_i is a 0 – 1-Bernoulli random variable and x_i is uniformly distributed over $[0, 1]^p$;
- for $i \in \mathcal{O}_4$, (x_i, y_i) is generated from the linear model (1) with the same true parameter $\beta = 1_p$ but for a different choice of design X and noise ε . Here, we take the covariance matrix Σ as an identity matrix and ε is a heavy-tailed noise following t_2 distribution.

After generating the corrupted data above, observations in \mathcal{I} and \mathcal{O} are merged and shuffled before downstream operations. Let the sample size $n = 5 \times 10^4$, the dimension $p = 20$, and the number of outliers $n_o = 19$. We subsample $r \in \{40p, 50p, 60p, 70p, 80p\}$ rows for row-sampling methods or $s = rp$ elements for MOM-CORE. The number of blocks is set to be $k = 40$ for the MOM procedure. Other settings are the same as those in the above section.

To evaluate the performance of different methods on corrupted data, we calculate MSE according to (10) and PMSE according to (11) for each method, and their results versus increasing subsample sizes are shown in Figs. 5 and 6, respectively. Remark that when predicting the responses, we first split the training set and test set on informative data $\{(x_i, y_i)\}_{i \in \mathcal{I}}$, and then add outliers $\{(x_i, y_i)\}_{i \in \mathcal{O}}$ to the training set; that is, the test set is not corrupted by outliers.

As shown in Figs. 5 and 6, MOM-CORE consistently achieves the smallest estimation and prediction errors among all these subsampling methods, and its advantage becomes more prominent with the increase of sparsity. Such an observation indicates that by integrating with the MOM procedure, our proposed MOM core-elements algorithm leads to an effective and robust estimator. Additionally, it is noteworthy that in both Figs. 3 and 6, the (MOM-)CORE prediction not only outperforms other subsampling approaches but also often delivers results that are nearly indistinguishable from full data.

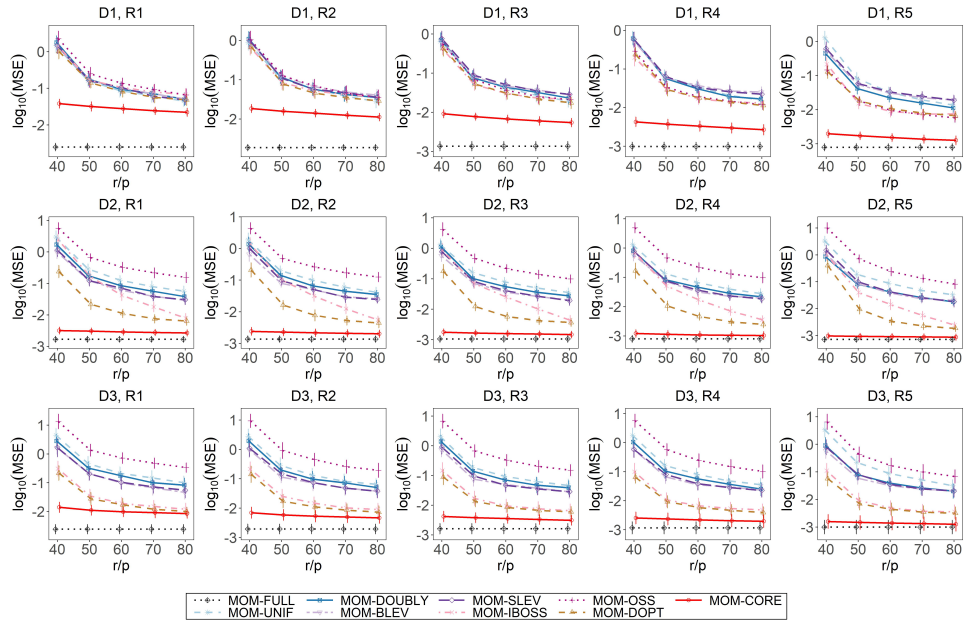


Figure 5: Comparison of different estimators w.r.t. MSE on corrupted data. Each row represents a particular data distribution, i.e., (D1)–(D3), and each column represents a different sparsity ratio, i.e., (R1)–(R5). Vertical bars are the standard errors.

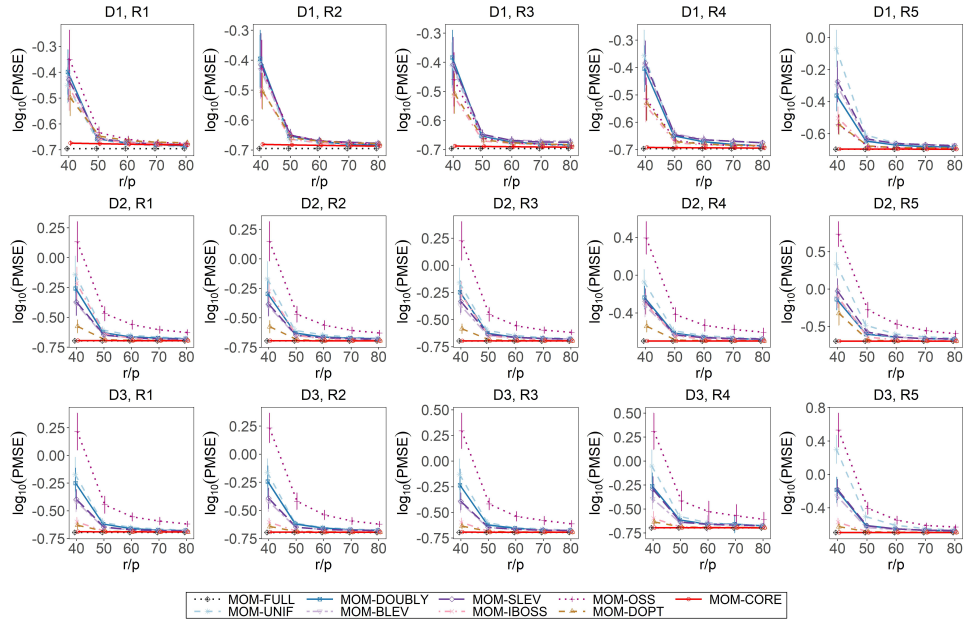


Figure 6: Comparison of different estimators w.r.t. PMSE on corrupted data. Each row represents a particular data distribution, i.e., (D1)–(D3), and each column represents a different sparsity ratio, i.e., (R1)–(R5). Vertical bars are the standard errors.

Table 2: CPU time (in seconds) of estimating β , for different combinations of n and p under (D1), (R1).

(a) CPU time for different p , with fixed $n = 10^5$.									
Method	FULL	UNIF	DOUBLY	BLEV	SLEV	IBOSS	OSS	DOPT	CORE
$p = 50$	0.82	0.00	0.02	0.51	0.51	0.15	0.40	0.43	0.08
$p = 100$	3.19	0.01	0.08	2.02	2.03	0.43	2.26	1.89	0.35
$p = 500$	76.31	0.39	1.61	46.37	46.39	6.53	58.77	45.54	5.29
(b) CPU time for different n , with fixed $p = 100$.									
Method	FULL	UNIF	DOUBLY	BLEV	SLEV	IBOSS	OSS	DOPT	CORE
$n = 5 \times 10^4$	1.72	0.00	0.06	1.01	1.01	0.12	1.87	0.89	0.13
$n = 5 \times 10^5$	16.95	0.02	0.10	10.86	10.85	1.36	6.41	9.92	0.87
$n = 5 \times 10^6$	189.65	0.11	0.21	119.83	119.84	12.94	43.28	91.52	7.54

5.3 Computing Time

To compare the computational efficiency of these subsampling approaches, we present the CPU time (in seconds) for different combinations of the sample size n and the dimension p under the case of (D1), (R1) in Table 2. Here, data corruption and the MOM procedure are omitted to save space, as they hardly affect the computation time. We take the subsample parameter $r = 10p$. All computations are implemented using the R programming language on a desktop running Windows 10 with an Intel i5-10210U CPU and 16GB memory. The CPU time for using the full sample is also presented for comparison.

As can be observed from Table 2, all of these estimates are more efficient than the full sample OLS estimate. It is unsurprising that the random sampling-based methods, UNIF and DOUBLY, require the least computing time, as they avoid the additional step of calculating subsampling probabilities. Among these methods, BLEV, SLEV, and DOPT necessitate the calculation of singular value decomposition of the full predictor matrix, leading to relatively longer CPU times. As expected, OSS also requires considerable computing time due to its complexity of $O(np \log r + rp^2)$.

Notably, the proposed core-elements approach only requires a longer time than UNIF and DOUBLY, while being faster than other competitors in almost all circumstances. In

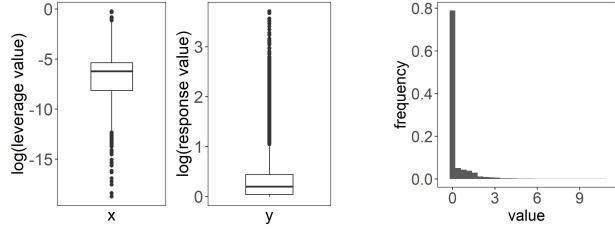
addition, its superiority in computation becomes more prominent when n is much larger than p , which is exactly the most suitable situation for taking advantage of subsampling.

6 Real Data Example

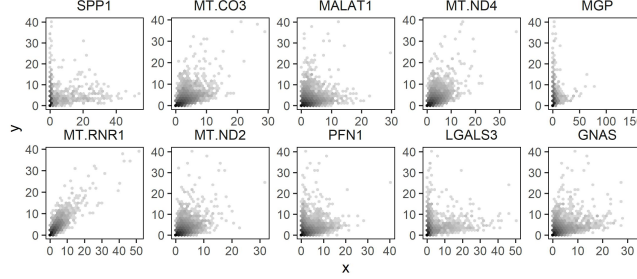
The rapid development of the single-cell RNA sequencing (scRNA-seq) technique enables the gene expression profiling of single cells. ScRNA-seq data are often organized into a reads count matrix, where rows are cells, columns represent genes, and the (i, j) th component is the observed expression level of the gene j in cell i . We consider a scRNA-seq dataset collected by [Azizi et al. \(2018\)](#), which includes CD 45+ immune cells from eight breast carcinomas, as well as matched normal breast tissue, blood, and lymph node. The dataset is publicly available with the accession code GSE114725 in Gene Expression Omnibus ([Edgar et al., 2002](#)). Our goal is to find the relationship between the expression of the gene MT-RNR2 and other genes. As a critical neuroprotective factor, the MT-RNR2 gene encodes the Humanin polypeptide and protects against death in Alzheimer’s disease.

To achieve the goal, we take the reads of this gene as the response and the reads for other genes as predictors. Following the data pre-processing steps in [Huang et al. \(2008\)](#), we first screen the genes as follows: (1) select the top 3000 genes with the highest expression levels; and (2) select the top 500 genes with the largest variances. We then standardize the predictors so that they have unit variance. The distribution of pre-processed data points is shown in Fig. 7(a). The final sample contains $n \approx 10^5$ cells and $p = 500$ genes with over 75% zero elements, as illustrated through the histogram in Fig. 7(b). Figure 7(c) presents the relationship between the response against ten randomly selected predictors of the scRNA-seq dataset, from which we can observe linear patterns. Therefore, it is reasonable to assume the data follow the linear model (1).

As the true coefficient vector β is unknown in real-world data analysis, we focus on comparing prediction performance by calculating the PMSE as defined in (11). This evaluation is based on one hundred bootstrap samples. The training and test sets are randomly partitioned according to the ratio of 7:3. The subsampling methods considered here are the same as those in Section 5, and the subsample size is set to be $r \in \{2p, 2^2p, 2^3p, 2^4p, 2^5p\}$ rows or $s = rp$ elements equivalently. We set the number of blocks for MOM to be $k = 5$.



(a) Box plots of predictor leverage values (left) and response values (right). (b) Histogram of values in the predictor matrix.



(c) Hexbin scatter plots (Carr et al., 2023) of the response against ten randomly selected predictors.

Figure 7: Visualization of the scRNA-seq dataset’s distribution. For visual clarity, the leverage and response values are (shifted to the positive range and) log-transformed in the subfigure (a), and the predictor values are truncated at the 99.9% quantile in the subfigure (b).

As shown in Fig. 7, the data exhibits linear patterns between the response and predictors, with no apparent outliers. Consequently, the MOM variant performs similarly to its original estimator. For clarity, we present only the MOM-FULL and MOM-CORE, omitting the MOM variants of other competing subsampling methods.

Figure 8 displays the performance of different approaches for predicting y on the test set. We observe the baseline methods, FULL and MOM-FULL, achieve almost the same prediction accuracy, suggesting the absence of extreme outliers in the dataset. Similarly, the proposed CORE and MOM-CORE methods have nearly identical prediction accuracy, both of which consistently outperform other subsampling approaches. In addition, the proposed estimators perform almost the same as $\hat{\beta}_{OLS}$ w.r.t. the PMSE, even when the selected number of elements is just $s = 2p^2$.

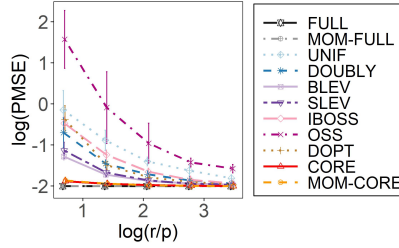


Figure 8: Comparison of different methods w.r.t. PMSE for the scRNA-seq dataset.

Table 3: CPU time (in seconds) of estimating β for the scRNA-seq dataset.

Method	FULL	UNIF	DOUBLY	BLEV	SLEV	IBOSS	OSS	DOPT	CORE
$r = 2p$	-	0.19	0.23	46.98	46.96	8.27	12.13	43.29	3.24
$r = 2^3p$	-	0.81	1.16	47.44	47.43	9.81	85.32	44.84	6.05
$r = 2^5p$	-	3.39	6.43	50.06	50.08	15.33	438.91	47.50	8.53
n	80.24	-	-	-	-	-	-	-	-

Table 3 shows the CPU time for these subsampling methods. Compared to the full sample estimate, all the subsample estimates require a shorter computational time, except for the OSS method as r increases. This exception can be attributed to the multiplicative order $O(\log r)$ in its computational cost. Similar to the results in Table 2, we observe that the core-elements approach has great advantages in computing. It is only slower than the random sampling-based UNIF and DOUBLY methods, while requiring nearly half the CPU time of the IBOSS method. The results in Fig. 8 and Table 3 indicate that the proposed strategy can provide a more effective estimate than the competitors, requiring a relatively short computational time.

7 Discussion

Realizing the gaps of element-wise subset selection methods in large-scale data analysis, we developed the core-elements method for least squares estimation in linear models. Theoretically, we showed that the proposed core-elements estimator approximately minimizes an upper bound of the estimation variance. We also provided a coresets-like finite-sample bound for the proposed estimator, which is supported by empirical results. To deal with

data corruption, we introduced the median-of-means estimation to provide a robust version of core-elements and established consistency of the resultant estimator. Empirical studies suggest that the proposed method is not only suitable for (numerically) sparse matrices but also has a superior performance for more general dense cases in terms of both accuracy and time.

Considering that the predictor matrix defined by basis function evaluations usually enjoys the (numerically) sparse property, we plan to extend core-elements to nonparametric additive models for efficiently approximating the penalized least squares estimation in the future. More importantly, the core-elements approach not only facilitates computation, but also has excellent applications in preserving data privacy and improving communication efficiency in federated learning, which are also left to our future work.

SUPPLEMENTARY MATERIAL

The Supplementary Material is structured as follows. Section 8 provides technical details of the theoretical results stated within the manuscript. The additional numerical results in Section 9 evaluate the performance of the proposed method in misspecified linear models and other choices of model coefficients.

8 Technical details

8.1 Proof of Lemma 1

Proof. The variance of $\tilde{\beta}$ can be organized as

$$\begin{aligned} E(\|\tilde{\beta} - \beta\|^2 | X) &= E[\{(X^{*\top} X)^{-1} X^{*\top} y - \beta\}^\top \{(X^{*\top} X)^{-1} X^{*\top} y - \beta\}] \\ &= E\{y^\top X^* (X^\top X^*)^{-1} (X^{*\top} X)^{-1} X^{*\top} y\} - \beta^\top \beta \\ &= \sigma^2 \|(X^{*\top} X)^{-1} X^{*\top}\|_F^2. \end{aligned} \tag{12}$$

The last equality is due to

$$E(y^\top A y) = (X\beta)^\top A X\beta + \sigma^2 \text{tr}(A)$$

with

$$A = X^* (X^\top X^*)^{-1} (X^{*\top} X)^{-1} X^{*\top}.$$

Recalling that $X^* = X - L$, the term $\|(X^{*\top}X)^{-1}X^{*\top}\|_F^2$ in (12) equals

$$\begin{aligned} & \text{tr}\{(X - L)(X^\top X - X^\top L)^{-1}(X^\top X - L^\top X)^{-1}(X - L)^\top\} \\ &= \text{tr}[(X - L)(X^\top X)^{-1}\{I_p - X^\top L(X^\top X)^{-1}\}^{-1}\{I_p - (X^\top X)^{-1}L^\top X\}^{-1}(X^\top X)^{-1}(X - L)^\top] \\ &= \text{tr}[\{I_p - X^\top L(X^\top X)^{-1}\}^{-1}\{I_p - (X^\top X)^{-1}L^\top X\}^{-1}(X^\top X)^{-1}(X - L)^\top(X - L)(X^\top X)^{-1}] \\ &= \text{tr}(S_1 S_2), \end{aligned}$$

where

$$S_1 = \{I_p - X^\top L(X^\top X)^{-1}\}^{-1}\{I_p - (X^\top X)^{-1}L^\top X\}^{-1}$$

and

$$S_2 = (X^\top X)^{-1}(X - L)^\top(X - L)(X^\top X)^{-1}.$$

Considering that S_2 is a positive semi-definite (PSD) matrix, it holds that

$$\text{tr}(S_1 S_2) \leq \sum_i \sigma_i(S_1) \sigma_i(S_2) \leq \sum_i \|S_1\|_2 \sigma_i(S_2) = \sum_i \|S_1\|_2 \lambda_i(S_2) = \|S_1\|_2 \text{tr}(S_2), \quad (13)$$

where $\sigma_i(\cdot)$ and $\lambda_i(\cdot)$ stand for the i th singular value and eigenvalue, respectively; the first inequality comes from the Von Neumann's trace inequality, and the last-but-one equality holds because the matrix S_2 is PSD. Therefore, it suffices to bound the terms $\|S_1\|_2$ and $\text{tr}(S_2)$.

First, we consider the term $\|S_1\|_2$. Performing a Taylor expansion of $\{I_p - (X^\top X)^{-1}L^\top X\}^{-1}$ around the point $L = 0_{n \times p}$ yields

$$\{I_p - (X^\top X)^{-1}L^\top X\}^{-1} = I_p + (X^\top X)^{-1}L^\top X + W_1$$

under the convergence condition that the spectral radius $\|(X^\top X)^{-1}L^\top X\|_2 = \lambda_0 < 1$; see [Higham \(2008, Chap. 1\)](#). Then, the remainder satisfies $\|W_1\|_2 = O(\lambda_0^2)$. Based on the result of Taylor expansion, we have

$$\begin{aligned} S_1 &= \{I_p + X^\top L(X^\top X)^{-1} + W_1^\top\} \{I_p + (X^\top X)^{-1}L^\top X + W_1\} \\ &= I_p + X^\top L(X^\top X)^{-1} + (X^\top X)^{-1}L^\top X + W_2 \end{aligned}$$

with

$$\begin{aligned} W_2 &= W_1 + W_1^\top + W_1^\top W_1 + W_1^\top (X^\top X)^{-1}L^\top X \\ &\quad + X^\top L(X^\top X)^{-1}W_1 + X^\top L(X^\top X)^{-2}L^\top X. \end{aligned}$$

One can see that $\|W_2\|_2 = O(\lambda_0^2)$. Consequently, we have

$$\|S_1\|_2 \leq 1 + O(\lambda_0). \quad (14)$$

Next, we consider the other term, $\text{tr}(S_2)$. It holds that

$$\text{tr}(S_2) = \text{tr}\{(X^\top X)^{-1}\} - 2 \text{tr}\{(X^\top X)^{-1} L^\top X (X^\top X)^{-1}\} + \text{tr}\{(X^\top X)^{-2} L^\top L\}. \quad (15)$$

As both $(X^\top X)^{-1}$ and $L^\top L$ are PSD matrices, we apply the inequality (13) on (15) and obtain that

$$\begin{aligned} \text{tr}(S_2) &\leq \text{tr}\{(X^\top X)^{-1}\} + 2\|(X^\top X)^{-1} L^\top X\|_2 \text{tr}\{(X^\top X)^{-1}\} + \|(X^\top X)^{-2}\|_2 \text{tr}(L^\top L) \\ &\leq \text{tr}\{(X^\top X)^{-1}\} + 2\lambda_0 \text{tr}\{(X^\top X)^{-1}\} + \|(X^\top X)^{-1}\|_2^2 \|L\|_F^2. \end{aligned} \quad (16)$$

By combining (14) and (16), we conclude that

$$\begin{aligned} \text{tr}(S_1 S_2) &\leq \|S_1\|_2 \text{tr}(S_2) \\ &\leq [\text{tr}\{(X^\top X)^{-1}\} + \|(X^\top X)^{-1}\|_2^2 \|L\|_F^2] \{1 + O(\lambda_0)\}. \end{aligned}$$

Therefore, the variance of $\tilde{\beta}$ can be bounded by

$$E(\|\tilde{\beta} - \beta\|^2 | X) \leq \sigma^2 [\text{tr}\{(X^\top X)^{-1}\} + \|(X^\top X)^{-1}\|_2^2 \|L\|_F^2] \{1 + O(\lambda_0)\}.$$

□

8.2 Proof of Theorem 2

Proof. (1) Recalling the definition of OLS estimation, i.e.,

$$\hat{\beta}_{\text{OLS}} = \arg \min_{\theta \in \mathbb{R}^p} \|y - X\theta\|^2,$$

the left-hand side of the inequality is obviously established.

(2) For the right-hand side of inequality,

$$\begin{aligned} \|y - X\tilde{\beta}\|^2 &= \|y - X\hat{\beta}_{\text{OLS}}\|^2 + 2(y - X\hat{\beta}_{\text{OLS}})^\top (X\hat{\beta}_{\text{OLS}} - X\tilde{\beta}) + \|X\hat{\beta}_{\text{OLS}} - X\tilde{\beta}\|^2 \\ &= \|y - X\hat{\beta}_{\text{OLS}}\|^2 + \|X\hat{\beta}_{\text{OLS}} - X\tilde{\beta}\|^2. \end{aligned}$$

The last equality holds because the cross term $(y - X\hat{\beta}_{\text{OLS}})^\top (X\hat{\beta}_{\text{OLS}} - X\tilde{\beta})$ equals to 0. Thus, it is sufficient to show that

$$\|X\hat{\beta}_{\text{OLS}} - X\tilde{\beta}\|^2 \leq \epsilon \|y - X\hat{\beta}_{\text{OLS}}\|^2. \quad (17)$$

Simple algebra yields that

$$\begin{aligned} \|X\hat{\beta}_{\text{OLS}} - X\tilde{\beta}\|^2 &= \|X(X^\top X)^{-1}X^\top y - X(X^{*\top}X)^{-1}X^{*\top}y\|^2 \\ &\leq \|X(X^\top X)^{-1}X^\top - X(X^{*\top}X)^{-1}X^{*\top}\|_2^2 \|y\|^2. \end{aligned} \quad (18)$$

Then, we consider the term $\|X(X^\top X)^{-1}X^\top - X(X^{*\top}X)^{-1}X^{*\top}\|_2^2$. Recall that $L = X - X^*$.

By using the condition $\|L\|_2 < \epsilon' \|X\|_2$, we have

$$\begin{aligned} &\|X(X^\top X)^{-1}X^\top - X(X^{*\top}X)^{-1}X^{*\top}\|_2 \\ &= \|X(X^\top X)^{-1}X^\top - X\{(X - L)^\top X\}^{-1}(X - L)^\top\|_2 \\ &= \|X\{(X^\top X)^{-1} - (X^\top X - L^\top X)^{-1}\}X^\top + X(X^\top X - L^\top X)^{-1}L^\top\|_2 \\ &\leq \|X\|_2^2 \|(X^\top X)^{-1} - (X^\top X - L^\top X)^{-1}\|_2 + \epsilon' \|X\|_2^2 \|(X^\top X - L^\top X)^{-1}\|_2 \\ &\leq \lambda_0 \|X\|_2^2 \|(X^\top X - L^\top X)^{-1}\|_2 + \epsilon' \|X\|_2^2 \|(X^\top X - L^\top X)^{-1}\|_2 \\ &= (\lambda_0 + \epsilon') \|X\|_2^2 \|(X^\top X - L^\top X)^{-1}\|_2, \end{aligned} \quad (19)$$

where $\lambda_0 = \|(X^\top X)^{-1}L^\top X\|_2$. The last inequality is due to $A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1}$ and then $\|A^{-1} - B^{-1}\|_2 \leq \|A^{-1}(B - A)\|_2 \|B^{-1}\|_2$.

Consider the term $\|(X^\top X - L^\top X)^{-1}\|_2$ in (19). A special case of Woodbury formula (Hager, 1989) takes the form $(A - B)^{-1} = A^{-1} + A^{-1}B(A - B)^{-1}$, which has a recursive structure that yields

$$(A - B)^{-1} = \sum_{k=0}^{\infty} (A^{-1}B)^k A^{-1}. \quad (20)$$

By using (20), it follows that

$$(X^\top X - L^\top X)^{-1} = \sum_{k=0}^{\infty} \{(X^\top X)^{-1}L^\top X\}^k (X^\top X)^{-1}.$$

Thus, we have

$$\begin{aligned} \|(X^\top X - L^\top X)^{-1}\|_2 &\leq \sum_{k=0}^{\infty} \|(X^\top X)^{-1}L^\top X\|_2^k \|(X^\top X)^{-1}\|_2 \\ &= \frac{1}{1 - \lambda_0} \|(X^\top X)^{-1}\|_2. \end{aligned}$$

In this way, the formula (19) can be further bounded by

$$\begin{aligned}\|X(X^\top X)^{-1}X^\top - X(X^{*\top}X)^{-1}X^{*\top}\|_2 &\leq \frac{\lambda_0 + \epsilon'}{1 - \lambda_0} \|X\|_2^2 \|(X^\top X)^{-1}\|_2 \\ &= \frac{\lambda_0 + \epsilon'}{1 - \lambda_0} \kappa^2(X),\end{aligned}\quad (21)$$

where $\kappa(X) = \sigma_{\max}(X)/\sigma_{\min}(X)$ is the condition number of X .

Considering the spectral radius λ_0 , we have

$$\lambda_0 = \|(X^\top X)^{-1}L^\top X\|_2 \leq \epsilon' \|(X^\top X)^{-1}\|_2 \|X\|_2^2 = \epsilon' \kappa^2(X). \quad (22)$$

Combining (18), (21) and (22), we conclude that if ϵ' satisfies

$$\frac{\{\kappa^2(X) + 1\}\epsilon'}{1 - \epsilon' \kappa^2(X)} \kappa^2(X) \leq \frac{\epsilon^{1/2} \|y - X\hat{\beta}_{\text{OLS}}\|}{\|y\|} \quad (23)$$

and

$$\epsilon' < \frac{1}{\kappa^2(X)}, \quad (24)$$

then the desired inequality (17) holds. The inequality (24) is required to ensure the upper bound (22) of λ_0 is smaller than 1. Further, the inequality (23) can be rewritten as

$$\epsilon' \leq \frac{1}{\kappa^2(X)} \left[1 + \frac{\{\kappa^2(X) + 1\} \|y\|}{\epsilon^{1/2} \|y - X\hat{\beta}_{\text{OLS}}\|} \right]^{-1} < \frac{1}{\kappa^2(X)}.$$

This completes the proof of the right-hand side of the inequality. \square

8.3 Proof of Remark 3

Proof. Recall that F is the cumulative distribution function of non-zero elements in X . Further assume the distribution is symmetric about zero, which is satisfied under the two specific cases we are concerned with.

We first consider the general continuous and symmetric F . Let $n_\alpha = \lfloor \alpha n \rfloor$. For $i = 1, \dots, n_\alpha$ and $j = 1, \dots, p$, define $z_{(i)j}$ be the i th order statistic for non-zero elements in $\{x_{1j}^2, \dots, x_{nj}^2\}$. As before, let $L = X - X^*$. Then, we have

$$\|L\|_2^2 \leq \|L\|_F^2 = \sum_{i=1}^{n_\alpha-r} \sum_{j=1}^p z_{(i)j} \leq (n_\alpha - r) \sum_{j=1}^p z_{(n_\alpha-r)j}.$$

Define $G(x) = F(\sqrt{x}) - F(-\sqrt{x})$ and $g(x) = dG(x)/dx$ for $x \geq 0$ as the cumulative distribution function and probability density function of non-zero elements in $\{x_{1j}^2, \dots, x_{nj}^2\}$,

respectively. By using the central limit theorem, $z_{(i)j}$ is asymptotically normal with mean $\mu_i = G^{-1}(i/n_\alpha)$ and variance $\sigma_i^2 = i(n_\alpha - i)/\{n_\alpha^3 g^2(\mu_i)\}$ (Walker, 1968). Thus, when $n_\alpha \rightarrow \infty$, we have the upper deviation inequality for $j = 1, \dots, p$,

$$\text{pr} \left(z_{(n_\alpha-r)j} \geq \mu_{n_\alpha-r} + t \right) \leq \exp \left(-\frac{t^2}{2\sigma_{n_\alpha-r}^2} \right) \quad \text{for any } t \geq 0,$$

that is,

$$\text{pr} \left[z_{(n_\alpha-r)j} < \mu_{n_\alpha-r} + \{-2\log(\delta)\}^{1/2} \sigma_{n_\alpha-r} \right] > 1 - \delta \quad \text{for any } 0 < \delta < 1.$$

Simple algebra yields that

$$\text{pr} \left(\|L\|_2^2 < (n_\alpha - r)p \left[\mu_{n_\alpha-r} + \{-2\log(\delta)\}^{1/2} \sigma_{n_\alpha-r} \right] \right) > (1 - \delta)^p.$$

Therefore, when the subsample parameter r in Algorithm 1 satisfies $r < \alpha n$ and

$$\frac{r}{n} \geq \alpha - \frac{\epsilon'^2 \|X\|_2^2}{np[\mu_{n_\alpha-r} + \{-2\log(\delta)\}^{1/2} \sigma_{n_\alpha-r}^2]}, \quad (25)$$

it holds that

$$\text{pr} \left(\|L\|_2^2 < \epsilon'^2 \|X\|_2^2 \right) > (1 - \delta)^p.$$

Consequently, $\|X - X^*\|_2 \leq \epsilon' \|X\|_2$ is achieved with probability at least $(1 - \delta)^p$ when n is sufficiently large, for any $0 < \delta < 1$.

Next, we focus on two specific cases.

- Uniform distribution over $(-1, 1)$, i.e., $F(x) = 2^{-1}(1 + x)\mathbb{I}_{(-1,1)}(x)$, where $\mathbb{I}(\cdot)$ is the indicator function.

Now we have $G(x) = \sqrt{x}\mathbb{I}_{(0,1)}(x)$ and $g(x) = (2\sqrt{x})^{-1}\mathbb{I}_{(0,1)}(x)$, which leads to $G^{-1}(x) = x^2\mathbb{I}_{(0,1)}(x)$. It follows that

$$\mu_{n_\alpha-r} = \left(\frac{n_\alpha - r}{n_\alpha} \right)^2 \quad \text{and} \quad \sigma_{n_\alpha-r}^2 = \frac{4(n_\alpha - r)^3 r}{n_\alpha^5}.$$

By plugging above results into (25), we conclude that

$$\frac{r}{n} \geq \alpha - \frac{(\alpha\epsilon'\|X\|_2)^{2/3}}{(2np)^{1/3}},$$

for a relatively large n .

- Standard normal distribution on \mathbb{R} .

Now the non-zero elements in $\{x_{1j}^2, \dots, x_{nj}^2\}$ follow the chi-squared distribution with 1 degree of freedom. Unfortunately, $\mu_{n_\alpha-r}$ and $\sigma_{n_\alpha-r}^2$ don't have closed-form expressions. According to the assumption $r > n_\alpha(1 - \phi)$, it holds that

$$\mu_{n_\alpha-r} < G^{-1}(\phi) \quad \text{and} \quad \sigma_{n_\alpha-r}^2 < \frac{(n_\alpha - r)r}{n_\alpha^3 g^2 \{G^{-1}(\phi)\}}.$$

Hence, when n is sufficiently large and

$$\frac{r}{n} \geq \alpha - \frac{(\epsilon' \|X\|_2)^2}{2G^{-1}(\phi)np},$$

the inequality (25) holds.

□

8.4 Proof of Theorem 4

Proof. Since the dimension of the parameter β is finite and the MOM procedure is adopted coordinate-wisely, without loss of generality, we only prove the univariate case, and we assume the first block contains no outliers.

By the definition of MOM, one can see that

$$\left\{ \|\tilde{\beta}_{\text{MOM}} - \beta\| > \epsilon \right\} \subseteq \left\{ \sum_{l=1}^k \mathbb{I}(\|\tilde{\beta}^{(l)} - \beta\| > \epsilon) \geq k/2 \right\},$$

where $\mathbb{I}(\cdot)$ is the indicator function and $\tilde{\beta}^{(l)}$ is the core-elements estimator on the l th block. To ease the conversation, denote $Z^{(l)} = \mathbb{I}(\|\tilde{\beta}^{(l)} - \beta\| > \epsilon)$ and $\mathcal{B}_{\mathcal{O}}^c$ be the complementary

set of $\mathcal{B}_\mathcal{O}$, i.e., $\mathcal{B}_\mathcal{O}^c = \{1, \dots, k\} \setminus \mathcal{B}_\mathcal{O}$. Then, we have

$$\begin{aligned}
& \text{pr}(\|\tilde{\beta}_{\text{MOM}} - \beta\| > \epsilon) \\
& \leq \text{pr} \left\{ \sum_{l=1}^k I(\|\tilde{\beta}^{(l)} - \beta\| > \epsilon) \geq k/2 \right\} \\
& = \text{pr} \left(\sum_{l \in \mathcal{B}_\mathcal{O}^c} Z^{(l)} + \sum_{l \in \mathcal{B}_\mathcal{O}} Z^{(l)} \geq k/2 \right) \\
& \leq \text{pr} \left(\sum_{l \in \mathcal{B}_\mathcal{O}^c} Z^{(l)} + |\mathcal{B}_\mathcal{O}| \geq k/2 \right) \\
& = \text{pr} \left[\frac{1}{k - |\mathcal{B}_\mathcal{O}|} \sum_{l \in \mathcal{B}_\mathcal{O}^c} \{Z^{(l)} - E(Z^{(l)})\} \geq \frac{k - 2|\mathcal{B}_\mathcal{O}|}{2(k - |\mathcal{B}_\mathcal{O}|)} - E(Z^{(1)}) \right] \tag{26}
\end{aligned}$$

$$\leq \exp \left[-2(k - |\mathcal{B}_\mathcal{O}|) \left\{ \frac{1}{2} - \frac{|\mathcal{B}_\mathcal{O}|}{2(k - |\mathcal{B}_\mathcal{O}|)} - E(Z^{(1)}) \right\}^2 \right], \tag{27}$$

where (26) comes from the fact that the partition is random, so that $\tilde{\beta}^{(l)}$ are i.i.d. for $l \in \mathcal{B}_\mathcal{O}^c$, which indicates $E(Z^{(l)}) = E(Z^{(1)})$ for all $l \in \mathcal{B}_\mathcal{O}^c$; the inequality (27) comes from the Hoeffding's inequality.

From Markov's inequality, one can see that

$$E(Z^{(1)}) \leq E(\|\tilde{\beta}^{(1)} - \beta\|^2)/\epsilon^2 \rightarrow 0$$

according to Lemma 1 under (H1)–(H3). Note that $|\mathcal{B}_\mathcal{O}|/(2k - 2|\mathcal{B}_\mathcal{O}|) \rightarrow 0$ since $|\mathcal{B}_\mathcal{O}|$ is finite and $k \rightarrow \infty$. Thus the desired result follows. \square

9 Additional numerical results

9.1 Misspecified linear model

In this section, we show the proposed core-elements estimator is robust to model misspecification. Suppose the true underlying model has the form

$$y_i = x_i^\top \beta + u_i, \quad i = 1, \dots, n. \tag{28}$$

Here, u_i 's are independently distributed random errors following the non-centered normal distribution $N(h(x_i), \sigma^2)$, where $h(\cdot)$ is an unknown multivariate function. Without prior

information on the true model (28), the classical linear model (1) in the manuscript is a misspecified linear model of (28).

We simulate the data from the model (28) with $n = 10^4$, $p = 20$ and $r \in \{2p, 4p, 6p, 8p, 10p\}$. The predictor matrix X is generated from (D1) with different sparsity patterns (R1)–(R5). Other settings are the same as those in subsection 5.1 of the manuscript. To show the robustness of the proposed method to various misspecification terms, following the work of Meng et al. (2021), we consider different kinds of $h(\cdot)$:

$$(M1) \quad h(x_i) = c_1 \cdot x_{i3}x_{i8};$$

$$(M2) \quad h(x_i) = c_2 \cdot x_{i3} \sin(x_{i8});$$

$$(M3) \quad h(x_i) = c_3 \cdot x_{i3}^2,$$

where the constants c_1, c_2 and c_3 are selected so that $\max_{x \in \{x_i\}_{i=1}^n} |h(x)| = 10$, that is, the misspecification term does not dominate the response.

Figures 9(a) and 9(b) compare the performance of different subsample estimators under model misspecification in estimation and prediction, respectively. We observe the proposed core-elements approach yields the most accurate estimation in all cases, indicating the performance of core-elements is robust to model misspecification.

9.2 Alternative choices of model coefficients

In subsampling literature, it is common practice to fix the true coefficient β as a constant vector, as seen in Wang et al. (2019, 2021); Yu et al. (2023); Chasiotis and Karlis (2024), and as employed in the main body of our manuscript. However, to comprehensively evaluate the adaptability of our proposed method in different scenarios, we also test it with alternative coefficients, including $\beta_i = (-1)^{i+1}, i = 1, \dots, p$ (Hou-Liu and Browne, 2023) and $\beta = (1_{10}, 0.1 \times 1_{p-20}, 1_{10})^\top$ (Ma et al., 2015). The remaining simulation settings are consistent with those in subsection 5.1.

The comparative results for the core-elements method and its competitors are presented in Figs. 10 and 11. In both cases, the performance of the full sample estimator and all subsampling methods show patterns highly similar to those observed in Figs. 2 and 3 of the manuscript, wherein our core-elements method not only surpasses other subsampling

approaches but also exhibits comparable efficacy to the full sample estimator, even at a relatively small subsample size (e.g., $r = 10p$). These findings highlight the generality and superiority of the core-elements approach across varied parametric landscapes.

References

- Achlioptas, D., Z. S. Karnin, and E. Liberty (2013). Near-optimal entrywise sampling for data matrices. *Advances in Neural Information Processing Systems* 26, 1565–1573.
- Achlioptas, D. and F. Mcsherry (2007). Fast computation of low-rank matrix approximations. *Journal of the Association for Computing Machinery* 54(2), 1–19.
- Ai, M., F. Wang, J. Yu, and H. Zhang (2020). Optimal subsampling for large-scale quantile regression. *Journal of Complexity* 62, 101512.
- Ai, M., J. Yu, H. Zhang, and H. Wang (2021). Optimal subsampling algorithms for big data regressions. *Statistica Sinica* 31, 749–772.
- Alaoui, A. and M. W. Mahoney (2015). Fast randomized kernel ridge regression with statistical guarantees. *Advances in Neural Information Processing Systems* 28, 775–783.
- Andrews, T. S., V. Y. Kiselev, D. McCarthy, and M. Hemberg (2021). Tutorial: guidelines for the computational analysis of single-cell RNA sequencing data. *Nature Protocols* 16(1), 1–9.
- Arora, S., E. Hazan, and S. Kale (2005). Fast algorithms for approximate semidefinite programming using the multiplicative weights update method. In *46th Annual IEEE Symposium on Foundations of Computer Science*, pp. 339–348. IEEE.
- Arora, S., E. Hazan, and S. Kale (2006). A fast random sampling algorithm for sparsifying matrices. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pp. 272–279. Springer.
- Azizi, E., A. J. Carr, G. Plitas, A. E. Cornish, C. Konopacki, S. Prabhakaran, J. Nainys, K. Wu, V. Kiseliovas, and M. Setty (2018). Single-cell map of diverse immune phenotypes in the breast tumor microenvironment. *Cell* 174(5), 1293–1308.

- Boutsidis, C., P. Drineas, and M. Magdon-Ismail (2013). Near-optimal coresets for least-squares regression. *IEEE Transactions on Information Theory* 59(10), 6880–6892.
- Braverman, V., R. Krauthgamer, A. R. Krishnan, and S. Sapir (2021). Near-optimal entrywise sampling of numerically sparse matrices. In *Proceedings of Thirty Fourth Conference on Learning Theory*, Volume 134, pp. 759–773. PMLR.
- Cai, J.-F., E. J. Candès, and Z. Shen (2010). A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization* 20(4), 1956–1982.
- Candès, E. J. and B. Recht (2009). Exact matrix completion via convex optimization. *Foundations of Computational Mathematics* 9(6), 717–772.
- Candès, E. J. and T. Tao (2010). The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory* 56(5), 2053–2080.
- Carmon, Y., Y. Jin, A. Sidford, and K. Tian (2020). Coordinate methods for matrix games. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science*, pp. 283–293. IEEE.
- Carr, D., N. Lewin-Koh, M. Maechler, and D. Sarkar (2023). hexbin: Hexagonal binning routines. *R package version 1.28.3*.
- Chang, M.-C. (2023). Predictive subdata selection for computer models. *Journal of Computational and Graphical Statistics* 32(2), 613–630.
- Chasiotis, V. and D. Karlis (2024). Subdata selection for big data regression: an improved approach. *Journal of Data Science, Statistics, and Visualisation* 4(3), 1–28.
- Chen, Y., S. Bhojanapalli, S. Sanghavi, and R. Ward (2014). Coherent matrix completion. In *International Conference on Machine Learning*, pp. 674–682. PMLR.
- Dai, W., Y. Song, and D. Wang (2023). A subsampling method for regression problems based on minimum energy criterion. *Technometrics* 65(2), 192–205.
- Dasgupta, A., P. Drineas, B. Harb, R. Kumar, and M. W. Mahoney (2009). Sampling algorithms and coresets for l_p regression. *SIAM Journal on Computing* 38(5), 2060–2078.

- Davis, T. A. and Y. Hu (2011). The University of Florida sparse matrix collection. *ACM Transactions on Mathematical Software* 38(1), 1–25.
- Dereziński, M., M. K. Warmuth, and D. J. Hsu (2018). Leveraged volume sampling for linear regression. *Advances in Neural Information Processing Systems* 31, 2510–2519.
- Donoho, D. L. and M. Gasko (1992). Breakdown properties of location estimates based on halfspace depth and projected outlyingness. *The Annals of Statistics* 20(4), 1803–1827.
- Donoho, D. L. and P. J. Huber (1983). The notion of breakdown point. *A Festschrift for Erich L. Lehmann*.
- Drineas, P., R. Kannan, and M. W. Mahoney (2006). Fast Monte Carlo algorithms for matrices I: Approximating matrix multiplication. *SIAM Journal on Computing* 36(1), 132–157.
- Drineas, P., M. Magdon-Ismail, M. W. Mahoney, and D. P. Woodruff (2012). Fast approximation of matrix coherence and statistical leverage. *The Journal of Machine Learning Research* 13(1), 3475–3506.
- Drineas, P. and A. Zouzias (2011). A note on element-wise matrix sparsification via a matrix-valued Bernstein inequality. *Information Processing Letters* 111(8), 385–389.
- d’Aspremont, A. (2011). Subsampling algorithms for semidefinite programming. *Stochastic Systems* 1(2), 274–305.
- Edgar, R., M. Domrachev, and A. E. Lash (2002). Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Research* 30(1), 207–210.
- El Karoui, N. and A. d’Aspremont (2010). Second order accurate distributed eigenvector computation for extremely large matrices. *Electronic Journal of Statistics* 4, 1345–1385.
- Feldman, D., M. Schmidt, and C. Sohler (2020). Turning big data into tiny data: Constant-size coresets for k -means, PCA, and projective clustering. *SIAM Journal on Computing* 49(3), 601–657.

- Garber, D. and E. Hazan (2016). Sublinear time algorithms for approximate semidefinite programming. *Mathematical Programming* 158(1), 329–361.
- Gupta, N. and A. Sidford (2018). Exploiting numerical sparsity for efficient learning: faster eigenvector computation and regression. *Advances in Neural Information Processing Systems* 31, 5274–5283.
- Hager, W. W. (1989). Updating the inverse of a matrix. *SIAM Review* 31(2), 221–239.
- Hampel, F. R. (1968). *Contributions to the Theory of Robust Estimation*. University of California, Berkeley.
- Hastie, T., R. Mazumder, J. D. Lee, and R. Zadeh (2015). Matrix completion and low-rank SVD via fast alternating least squares. *Journal of Machine Learning Research* 16(1), 3367–3402.
- Higham, N. J. (2008). *Functions of Matrices: Theory and Computation*. SIAM.
- Hou-Liu, J. and R. P. Browne (2023). Generalized linear models for massive data via doubly-sketching. *Statistics and Computing* 33(105), 1–19.
- Hsu, D. and S. Sabato (2014). Heavy-tailed regression with a generalized median-of-means. In *International Conference on Machine Learning*, pp. 37–45. PMLR.
- Hu, Y., M. Li, X. Liu, and C. Meng (2024). Sampling-based methods for multi-block optimization problems over transport polytopes. *Mathematics of Computation*, In press.
- Huang, J., S. Ma, and C. H. Zhang (2008). Adaptive Lasso for sparse high-dimensional regression models. *Statistica Sinica* 18(4), 1603–1618.
- Huang, S.-T. and J. Lederer (2023). Deepmom: Robust deep learning with median-of-means. *Journal of Computational and Graphical Statistics* 32(1), 181–195.
- Joseph, V. R. and S. Mak (2021). Supervised compression of big data. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 14(3), 217–229.
- Joseph, V. R. and A. Vakayil (2022). SPlit: An optimal method for data splitting. *Technometrics* 64(2), 166–176.

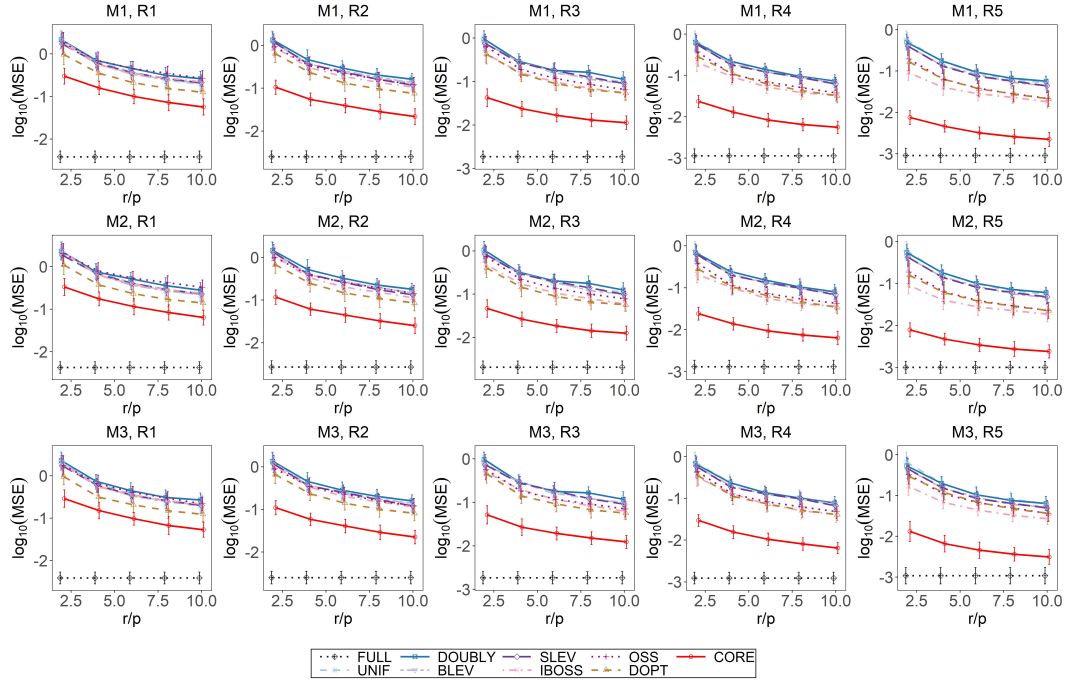
- Kairouz, P., H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, et al. (2021). Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning* 14(1–2), 1–210.
- Knight, K. (2018). Subsampling least squares and elemental estimation. In *2018 IEEE Data Science Workshop (DSW)*, pp. 91–94. IEEE.
- Konečný, J., H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon (2016). Federated learning: Strategies for improving communication efficiency. In *NIPS Workshop on Private Multi-Party Machine Learning*.
- Kundu, A., P. Drineas, and M. Magdon-Ismail (2017). Recovering PCA and sparse PCA via hybrid- (ℓ_1, ℓ_2) sparse sampling of data elements. *Journal of Machine Learning Research* 18(1), 2558–2591.
- Lecué, G. and M. Lerasle (2019). Learning from MOM’s principles: Le Cam’s approach. *Stochastic Processes and Their Applications* 129(11), 4385–4410.
- Lecué, G. and M. Lerasle (2020). Robust machine learning by median-of-means: theory and practice. *The Annals of Statistics* 48(2), 906–931.
- Li, F., R. Xie, Z. Wang, L. Guo, J. Ye, P. Ma, and W. Song (2019). Online distributed IoT security monitoring with multidimensional streaming big data. *IEEE Internet of Things Journal* 7(5), 4387–4394.
- Li, M., J. Yu, T. Li, and C. Meng (2023a). Importance sparsification for Sinkhorn algorithm. *Journal of Machine Learning Research* 24, 1–44.
- Li, M., J. Yu, H. Xu, and C. Meng (2023b). Efficient approximation of Gromov-Wasserstein distance using importance sparsification. *Journal of Computational and Graphical Statistics* 32(4), 1512–1523.
- Li, M., J. Zhang, and C. Meng (2024). Nonparametric additive models for billion observations. *Journal of Computational and Graphical Statistics*, In press.
- Li, T. and C. Meng (2021). Modern subsampling methods for large-scale least squares regression. *International Journal of Cyber-Physical Systems* 2(2), 1–28.

- Liu, F., B. Zhu, S. Yuan, J. Li, and K. Xue (2021). Privacy-preserving truth discovery for sparse data in mobile crowdsensing systems. In *2021 IEEE Global Communications Conference*, pp. 1–6. IEEE.
- Lugosi, G. and S. Mendelson (2019). Regularization, sparse recovery, and median-of-means tournaments. *Bernoulli* 25(3), 2075–2106.
- Ma, P., J. Z. Huang, and N. Zhang (2015). Efficient computation of smoothing splines via adaptive basis sampling. *Biometrika* 102(3), 631–645.
- Ma, P., M. W. Mahoney, and B. Yu (2015). A statistical perspective on algorithmic leveraging. *Journal of Machine Learning Research* 16(1), 861–911.
- Ma, P. and X. Sun (2015). Leveraging for big data regression. *Wiley Interdisciplinary Reviews: Computational Statistics* 7(1), 70–76.
- Ma, P., X. Zhang, X. Xing, J. Ma, and M. Mahoney (2020). Asymptotic analysis of sampling estimators for randomized numerical linear algebra algorithms. In *International Conference on Artificial Intelligence and Statistics*, pp. 1026–1035. PMLR.
- Maalouf, A., G. Eini, B. Mussay, D. Feldman, and M. Osadchy (2022). A unified approach to coresets learning. *IEEE Transactions on Neural Networks and Learning Systems*, 1–13.
- Mak, S. and V. R. Joseph (2018). Support points. *The Annals of Statistics* 46(6A), 2562–2592.
- Martinez, C. (2004). Partial quicksort. In *Proc. 6th ACM-SIAM Workshop on Algorithm Engineering and Experiments and 1st ACM-SIAM Workshop on Analytic Algorithmics and Combinatorics*, pp. 224–228.
- Mathieu, T. (2021). *M-estimation and Median of Means applied to statistical learning*. Ph.D. thesis, Université Paris-Saclay.
- Meng, C., Y. Wang, X. Zhang, A. Mandal, W. Zhong, and P. Ma (2017). Effective statistical methods for big data analytics. In *Handbook of Research on Applied Cybernetics and Systems Science*, pp. 280–299. IGI Global.

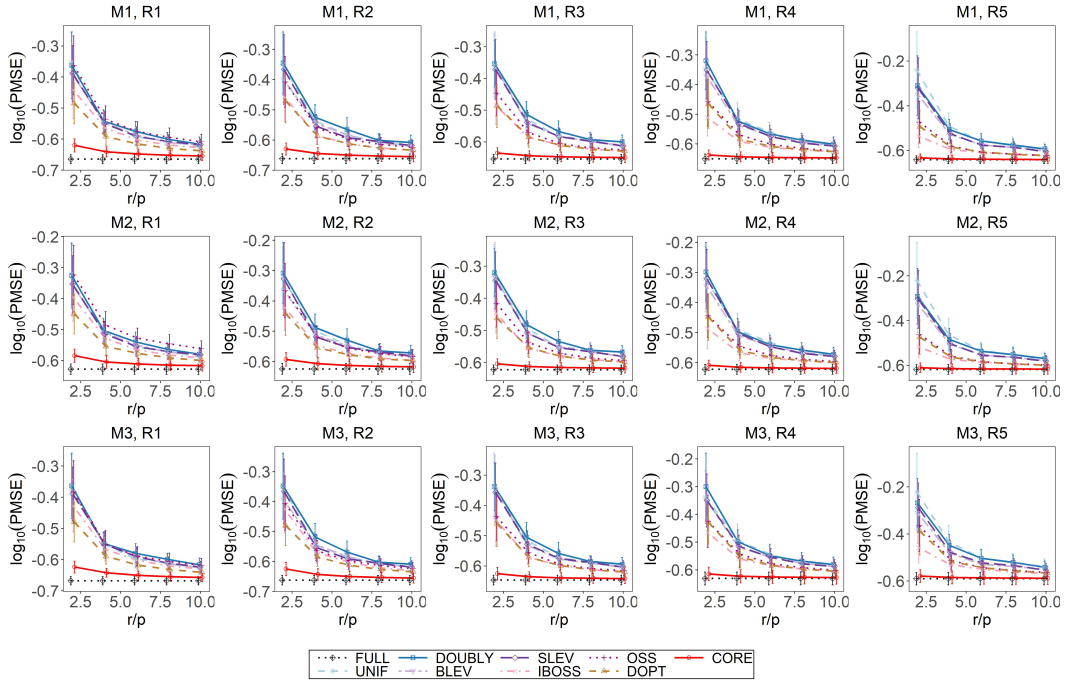
- Meng, C., R. Xie, A. Mandal, X. Zhang, W. Zhong, and P. Ma (2021). Lowcon: A design-based subsampling approach in a misspecified linear model. *Journal of Computational and Graphical Statistics* 30(3), 694–708.
- Meng, C., J. Yu, Y. Chen, W. Zhong, and P. Ma (2022). Smoothing splines approximation using Hilbert curve basis selection. *Journal of Computational and Graphical Statistics*, 1–11.
- Meng, C., X. Zhang, J. Zhang, W. Zhong, and P. Ma (2020). More efficient approximation of smoothing splines via space-filling basis selection. *Biometrika* 107, 723–735.
- Munteanu, A., C. Schwiiegelshohn, C. Sohler, and D. P. Woodruff (2018). On coresets for logistic regression. *Advances in Neural Information Processing Systems* 31, 6562–6571.
- Musser, D. R. (1997). Introspective sorting and selection algorithms. *Software: Practice and Experience* 27(8), 983–993.
- Muzellec, B., J. Josse, C. Boyer, and M. Cuturi (2020). Missing data imputation using optimal transport. In *International Conference on Machine Learning*, pp. 7130–7140. PMLR.
- Nguyen, T. K. H., K. Van den Berge, M. Chiogna, and D. Risso (2023). Structure learning for zero-inflated counts with an application to single-cell RNA sequencing data. *The Annals of Applied Statistics* 17(3), 2555–2573.
- Qaiser, S. and R. Ali (2018). Text mining: Use of TF-IDF to examine the relevance of words to documents. *International Journal of Computer Applications* 181(1), 25–29.
- Reuter, T. and R. Schwabe (2024). D-optimal subsampling design for massive data linear regression. *arXiv preprint arXiv:2307.02236*.
- Settles, B. (2012). Active learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning* 6(1), 1–114.
- Sun, X., W. Zhong, and P. Ma (2021). An asymptotic and empirical smoothing parameters selection method for smoothing spline anova models in large samples. *Biometrika* 108(1), 149–166.

- Vakayil, A. and V. R. Joseph (2022). Data twinning. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 15(5), 598–610.
- Van Buuren, S. and K. Groothuis-Oudshoorn (2011). MICE: Multivariate imputation by chained equations in R. *Journal of Statistical Software* 45, 1–67.
- Walker, A. (1968). A note on the asymptotic distribution of sample quantiles. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 30(3), 570–575.
- Wang, H., M. Yang, and J. Stufken (2019). Information-based optimal subdata selection for big data linear regression. *Journal of the American Statistical Association* 114(525), 393–405.
- Wang, H., R. Zhu, and P. Ma (2018). Optimal subsampling for large sample logistic regression. *Journal of the American Statistical Association* 113(522), 829–844.
- Wang, J., J. Zou, and H. Wang (2022). Sampling with replacement vs Poisson sampling: a comparative study in optimal subsampling. *IEEE Transactions on Information Theory* 68(10), 6605–6630.
- Wang, L., J. Elmstedt, W. K. Wong, and H. Xu (2021). Orthogonal subsampling for big data linear regression. *The Annals of Applied Statistics* 15(3), 1273–1290.
- Wang, R., Y. Ouyang, Y. Panpan, and W. Xu (2023). A fast and accurate estimator for large scale linear model via data averaging. In *Advances in Neural Information Processing Systems*, Volume 36, pp. 34917–34927.
- Wang, S., A. Gittens, and M. W. Mahoney (2019). Scalable kernel K-means clustering with Nyström approximation: relative-error bounds. *Journal of Machine Learning Research* 20(1), 431–479.
- Wang, S. and Z. Zhang (2013). Improving CUR matrix decomposition and the Nyström approximation via adaptive sampling. *Journal of Machine Learning Research* 14(1), 2729–2769.

- Xie, R., Z. Wang, S. Bai, P. Ma, and W. Zhong (2019). Online decentralized leverage score sampling for streaming multidimensional time series. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2301–2311.
- Yu, J., M. Ai, and Z. Ye (2023). A review on design inspired subsampling for big data. *Statistical Papers*, 1–44.
- Yu, J., H. Wang, M. Ai, and H. Zhang (2022). Optimal distributed subsampling for maximum quasi-likelihood estimators with massive data. *Journal of the American Statistical Association* 117(537), 265–276.

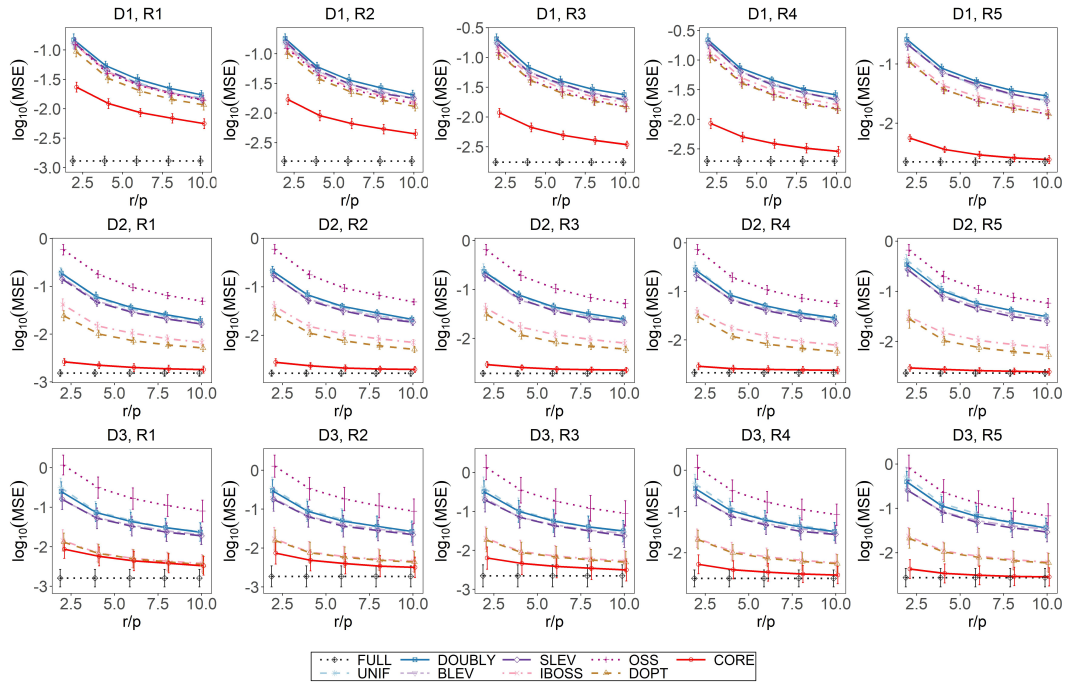


(a) Estimation performance.

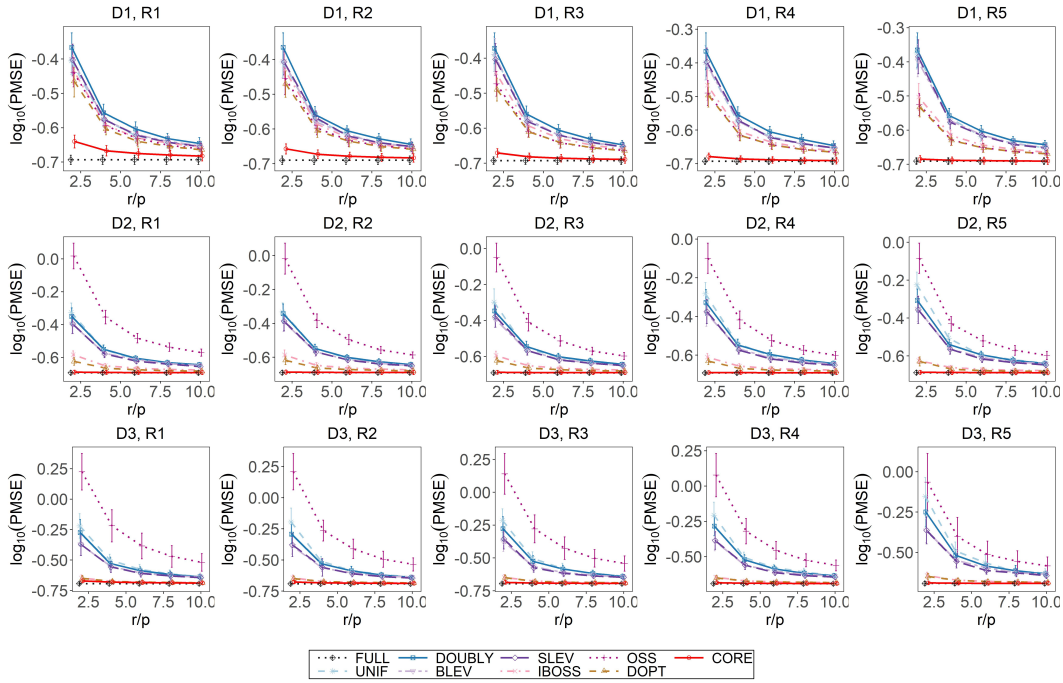


(b) Prediction performance.

Figure 9: Comparison of different estimators w.r.t. MSE and PMSE. Each row represents a particular misspecification term, i.e., (M1)–(M3), and each column represents a different sparsity ratio, i.e., (R1)–(R5). Vertical bars are the standard errors.

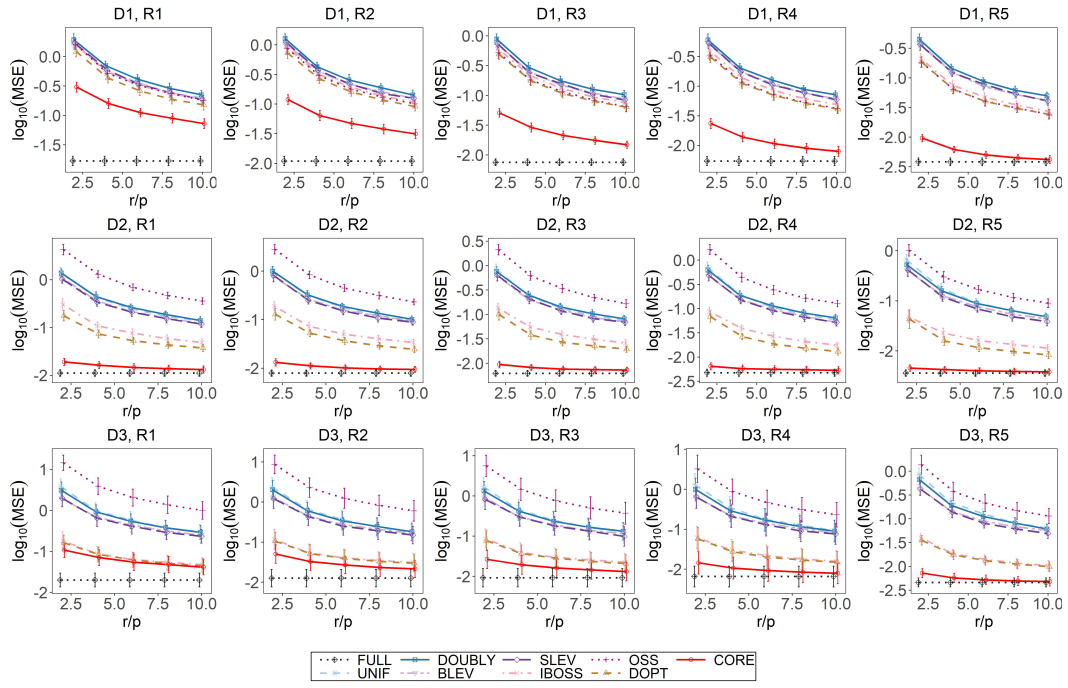


(a) Estimation performance.

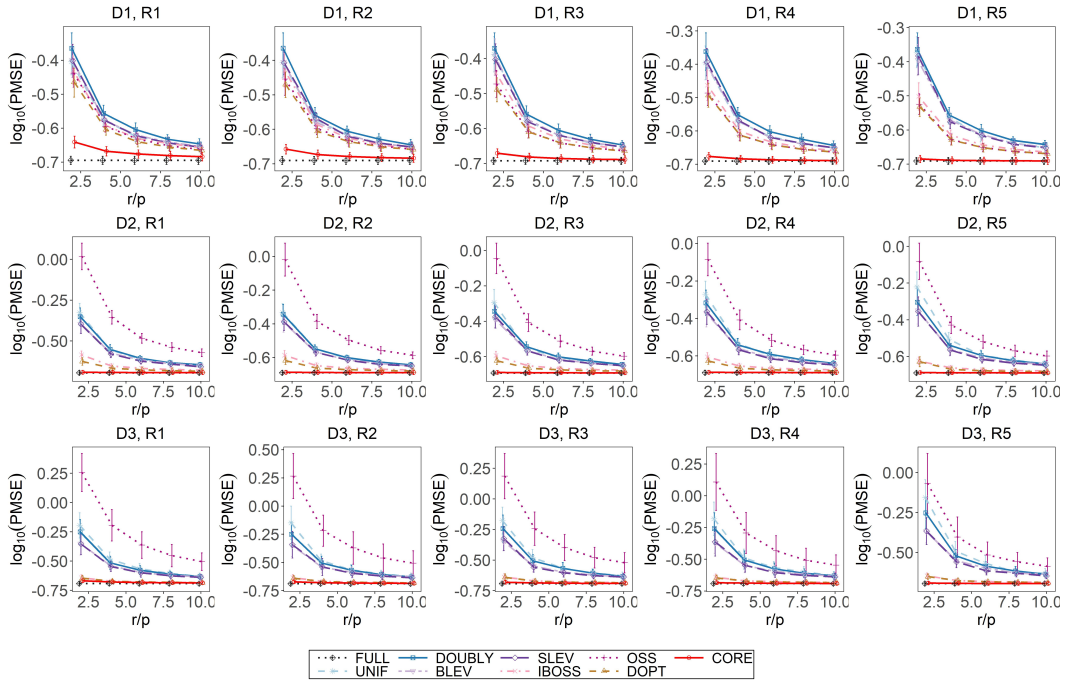


(b) Prediction performance.

Figure 10: Comparison of different estimators w.r.t. MSE and PMSE for $\beta_i = (-1)^{i+1}, i = 1, \dots, p$. Each row represents a particular data distribution, i.e., (D1)–(D3), and each column represents a different sparsity ratio, i.e., (R1)–(R5). Vertical bars are the standard errors.



(a) Estimation performance.



(b) Prediction performance.

Figure 11: Comparison of different estimators w.r.t. MSE and PMSE for $\beta = (1_{10}, 0.1 \times 1_{p-20}, 1_{10})^\top$. Each row represents a particular data distribution, i.e., (D1)–(D3), and each column represents a different sparsity ratio, i.e., (R1)–(R5). Vertical bars are the standard errors.